



Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT.

Theory in Biosciences 123 (2005) 301–369

Theory in  
Biosciences

[www.elsevier.de/thbio](http://www.elsevier.de/thbio)

## MATH/CHEM/COMP 2004

### Evolutionary patterns of non-coding RNAs

Athanasius F. Bompfünnewerer<sup>a,b</sup>, Christoph Flamm<sup>a</sup>,  
Claudia Fried<sup>c</sup>, Guido Fritzsch<sup>d</sup>, Ivo L. Hofacker<sup>a</sup>,  
Jörg Lehmann<sup>c</sup>, Kristin Missal<sup>c</sup>, Axel Mosig<sup>c</sup>, Bettina Müller<sup>c,e</sup>,  
Sonja J. Prohaska<sup>c</sup>, Bärbel M.R. Stadler<sup>f</sup>, Peter F. Stadler<sup>a,c,d,g,\*</sup>,  
Andrea Tanzer<sup>a,d</sup>, Stefan Washietl<sup>a</sup>, Christina Witwer<sup>a</sup>

<sup>a</sup>Department of Theoretical Chemistry, University of Vienna, Währingerstraße 17,  
A-1090 Wien, Austria

<sup>b</sup>Zentralfriedhof Wien, 3. Tor, Simmeringer Haupstraße 234, A-1110 Wien, Austria

<sup>c</sup>Bioinformatics Group, Department of Computer Science, University of Leipzig,  
Härtelstraße 16-18, D-04107 Leipzig, Germany

<sup>d</sup>Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstraße 16-18,  
D-04107 Leipzig, Germany

<sup>e</sup>Department of Biotechnology and Bioinformatics, University of Applied Sciences,  
Weihenstephan, D-85350 Freising, Germany

<sup>f</sup>Max Planck Institute for Mathematics in the Sciences, Inselstrasse 22-26, D-04103 Leipzig,  
Germany

<sup>g</sup>Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, NM 87501, USA

Received 22 December 2004; accepted 24 January 2005

### Abstract

A plethora of new functions of non-coding RNAs (ncRNAs) have been discovered in past few years. In fact, RNA is emerging as the central player in cellular regulation, taking on active roles in multiple regulatory layers from transcription, RNA maturation, and RNA

\*Corresponding author. Bioinformatics Group, Department of Computer Science, and Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstraße 7b, D-04107 Leipzig, Germany. Tel.: +49 341 97 16691; fax: +49 341 97 16709.

E-mail address: peter.stadler@bioinf.uni-leipzig.de (P.F. Stadler).

modification to translational regulation. Nevertheless, very little is known about the evolution of this “Modern RNA World” and its components. In this contribution, we attempt to provide at least a cursory overview of the diversity of ncRNAs and functional RNA motifs in non-translated regions of regular messenger RNAs (mRNAs) with an emphasis on evolutionary questions. This survey is complemented by an in-depth analysis of examples from different classes of RNAs focusing mostly on their evolution in the vertebrate lineage. We present a survey of Y RNA genes in vertebrates and study the molecular evolution of the U7 snRNA, the snoRNAs E1/U17, E2, and E3, the Y RNA family, the *let-7* microRNA (miRNA) family, and the mRNA-like *eif-1* gene. We furthermore discuss the statistical distribution of miRNAs in metazoans, which suggests an explosive increase in the miRNA repertoire in vertebrates. The analysis of the transcription of ncRNAs suggests that small RNAs in general are genetically mobile in the sense that their association with a hostgene (e.g. when transcribed from introns of a mRNA) can change on evolutionary time scales. The *let-7* family demonstrates, that even the mode of transcription (as intron or as exon) can change among paralogous ncRNA.

© 2005 Elsevier GmbH. All rights reserved.

**Keywords:** Evolution; Non-coding RNA; mRNA; rRNA; snRNA; snoRNA; miRNA; Y-RNA; Vault RNA; gRNA; RNA editing; UTR

---

## Introduction

Although it is still commonplace to speak of “*genes and their encoded protein products*”, thousands of human genes produce transcripts that exert their function without ever producing proteins. The diversity of sequences, sizes, structures, and functions of the known non-coding RNAs (ncRNAs) strongly suggests that we have seen only a small fraction of the functional RNAs. Most of the ncRNAs are small, they do not have translated ORFs, and they are not polyadenylated. Unlike protein coding genes, ncRNA gene sequences do not exhibit a strong *common* statistical signal, hence a reliable general purpose computational genefinder for ncRNA genes has been elusive (Eddy, 2001).

The list of functional ncRNAs includes key players in the biochemistry of the cell. Many of them have characteristic secondary structures that are highly conserved in evolution. A non-exhaustive list is compiled in Table 1. In addition to these relatively well-described classes there is a diverse and rapidly growing list of ncRNAs with sometimes enigmatic function: the 17 kb *Xist* RNA of humans and the smaller *roX* RNAs of *Drosophila* play a key role in dosage compensation and X chromosome inactivation (Avner and Heard, 2001; Franke and Baker, 2000). Several large ncRNAs are expressed from imprinted regions, see also Seitz et al. (2004). Many of these are *cis*-antisense RNAs that overlap coding genes on the other genomic strand (Erdmann et al., 2001). An RNA (meiRNA) regulates the onset of meiosis in fission yeast (Ohno and Mattaj, 1999). No precise function is known at present for the human H19 transcript, or the hrs $\omega$  transcript induced by heat shock in *Drosophila*, see e.g. Erdmann et al. (1999). A recent survey of the slime mold *Dictyostelium discoideum* uncovered two novel classes of ncRNAs (Aspegren et al., 2004). An

**Table 1.** Major classes of functional RNAs

Class	Size	Function	Phylogenetic distribution	DB
tRNA	70–80	Translation	Ubiquitous	Sprinzl et al. (1998)
rRNA	16S/18S	1.5k	Translation	Van de Peer et al. (2000b) and Maidak et al. (2001)
	28S+5.8S/23S	3k	Translation	Maidak et al. (2001)
RNase P	5S	130	Translation	Wuyts et al. (2001) and Maidak et al. (2001)
	P	220–440	tRNA maturation	Szymański et al. (2000)
	MRP	250–350	Endonuclease, 5.8S rRNA maturation	Brown (1999)
snoRNA	H/ACA	~130	Pseudouridinylation in rRNAs	Samarsky and Fournier (1999)
	C/D	60–80	Ribose 2'-O-methylation in rRNAs	Eukarya, archaea
snRNA	Telomerase	400–550	Major spliceosome, mRNA maturation	Eukarya
	U1,U2,U4,U5,U6	100–160	Minor spliceosome, mRNA maturation	Eukarya
	U11,U12	130–140	<i>trans</i> -splicing	Eukarya
	SL	~100	Histone mRNA maturation	Lower eukaryotes
	U7	~65	Transcriptional regulation	Eukarya
	7SK	~300	Signal recognition particle	Vertebrata
	7SL/SRP	300–400	Part of vault particle	Ubiquitous
	Vault	80–100	Part of Ro particle	Vertebrata
	Y	80–100	Tags protein for proteolysis	Metazoa
	300–400		Bacteria, chloroplasts, cyanoplasts	Zwieb and Wower (2000)
tmRNA	~22		Post-transcriptional regulation	Multicellular organisms
miRNA			RNA editing	Griffiths-Jones et al. (2003)
gRNA	40–80		Kineto/plastids	Hinz and Göringer (1999)

experimental screen recovered hundreds of small ncRNAs from the mouse (Hüttenhofer et al., 2001). Ambros et al. (2003) reported more than 30 tiny ncRNAs in a recent survey of *Caenorhabditis elegans* that are slightly shorter than microRNAs (miRNAs), are not processed from hairpin precursors, and are poorly conserved between related species.

Since the discovery of miRNAs (Lagos-Quintana et al., 2001; Lau et al., 2001; Lee and Ambros, 2001) and the development of RNAi as a general technique for manipulating translation (Elbashir et al., 2001), there is mounting evidence that ncRNAs in fact dominate the regulatory networks of the cell (Bartel and Chen, 2004; Hobert, 2004; Mattick, 2003, 2004; Szymbański et al., 2003): The *Escherichia coli* genome encodes more than 50 small RNA genes at least some of which (e.g., *MicF*, *OxyS*, *DsrA*, *Spot42*, *RhyB*) act by base-pairing to activate or repress translation (Gottesman, 2004; Storz et al., 2004). A large fraction of the mouse transcriptome consists of ncRNAs, many of them anti-sense to known protein-coding transcripts (Suzuki and Hayashizaki, 2004). Similarly, about half of the transcripts from Human chromosomes 21 and 22 are non-coding (Cawley et al., 2004; Kampa et al., 2004), see Morey and Avner (2004) for a discussion of the possible roles of anti-sense RNAs. *Leishmania* and related kinetoplastids have reduced transcriptional regulation of gene expression to a minimum, maybe to the point of having lost any specific polymerase II transcription initiation (Clayton, 2002). Instead, these organisms use an elaborate cleavage and *trans*-splicing mechanism based on the action of ~40 nt “spliced leader” RNA. *Tetrahymena* appears to use an RNA-based mechanism for directing its genome-wide DNA rearrangements (Mochizuki et al., 2002; Yao et al., 2003).

Another level of RNA function is presented by functional motifs within protein-coding RNAs. We briefly mention a few of the best-understood examples of structurally conserved RNA motifs in viral RNAs: an internal ribosomal entry site (IRES) region is used instead of a cap to initialize translation by Picornaviridae, some Flaviviridae including Hepatitis C virus, and a small number of messenger RNAs (mRNAs), see e.g. Rueckert (1996), Huez et al. (1998) and Pesole et al. (2001). Viral RNAs contain a large number of structured binding motifs that are essential for the viral life cycles, e.g., the TAR and RRE motifs in HIV (Dayton et al., 1992) or the *cis*-acting replication element (CRE) hairpin in Picornaviridae (Witwer et al., 2001). RNA-localization mechanisms involve specific sequences motifs in the localized RNA that cause certain proteins to mediate the interaction with cytoskeletal elements (Oleynikov and Singer, 1998). The localized *bicoid* mRNA, for instance, is responsible for laying down the body axes of the embryo (Pokrywka and Stephenson, 1991).

*RNA switches*, i.e., RNAs that drastically change their structure, are important regulatory elements (Sullenger, 2004). For instance, the terminator and anti-terminator, two alternative RNA hairpins, regulate gene expression in *E. coli* and *Bacillus subtilis* by attenuation (Babitzke and Yanofsky, 1993; Fayat et al., 1983; Putzer et al., 1992). RNA switches can provide exact temporal control as in the *hok/sok* system of plasmid R1 which triggers programmed cell death (Nagel et al., 1999; Møller-Jensen et al., 2001). RNA switches also play a role in the spliced leader of

trypanosomes and nematodes (LeCuyer et al., 1994). A theoretical study shows that RNAs exhibiting very different secondary structures with near-groundstate energy, i.e., potential riboswitches, are relatively frequent and easily accessible in evolution (Flamm et al., 2000). Artificial riboswitches have been explored for biotechnological applications (Soukup and Breaker, 1999; Komatsu, 2004; Schultes and Bartel, 2000) and it has been demonstrated that such constructs can be specifically triggered by means of small “modifier” RNAs (Meisner et al., 2004; Hackermüller et al., 2005).

Given the importance of ncRNAs and RNA-based mechanism in extant lifeforms, it is surprising that we know relatively little about the evolutionary history of most RNA classes. There are strong reasons to conclude that the Last Common Ancestor (LCA) was preceded by simpler life forms that were based primarily on RNA. In this *RNA World* scenario (Gilbert, 1986; Gesteland and Atkins, 1993), the translation of RNA into proteins and, finally, the usage of DNA (Freeland et al., 1999) as information storage device are later innovations. The wide range of catalytic activities that can be realized by relatively small ribozymes (Baskerville and Bartel, 2002; Illangasekare and Yarus, 1999; Johnston et al., 2001; Joyce, 2002; Lee et al., 2000; Unrau and Bartel, 1998) as well as the usage of RNA catalysis at crucial points of the information metabolism of modern cells (Jeffares et al., 1998; Doudna and Cech, 2002; Moore and Steitz, 2002) provides support for the *RNA World* hypothesis. Plausible ribozyme catalyzed pathways for a late-stage ribo-organism (Joyce, 2002), the role and evolution of co-enzymes (Jadhav and Yarus, 2002), and a rather detailed model of the steps leading from the RNA world to modern cellular architectures (Poole et al., 2000) have been the subject of detailed investigations.

Probably the best-studied group are the ribosomal RNAs (rRNAs) because of their utility in molecular phylogenetics. In fact, much of our knowledge about the deepest branches of the tree of life has been inferred from 16S/18S sequence data (Doolittle and Brown, 1994; Olsen and Woese, 1993; Van de Peer et al., 2000a; Peterson and Eernisse, 2001; Cavalier-Smith and Chao, 2003). Besides the 16S/18S and the 28S/23S large subunit rRNA, other classes of RNAs, however, have been used only sporadically for these purposes, although it has been shown that they are phylogenetically informative (Caetano-Anollés, 2002a; Collins et al., 2000; Hudelot et al., 2003). Telomerase RNA structures were used to elucidate the phylogeny of tetrahymenine ciliates (Ye and Romero, 2002). Nevertheless, relatively little information is available on the origins of various RNA classes. Apart from the rRNAs (see e.g. Caetano-Anollés (2002b)) and transfer RNAs (tRNAs) (Eigen et al., 1989), an origin predating the last common ancestor is clear only for the RNase P/RNase MRP family.

The Rfam data base (Griffiths-Jones et al., 2003, 2005), the non-code data base (Liu et al., 2005), and the RNAdb (Pang et al., 2005) collect the flood of information on such ncRNAs and functional RNA motifs that before has been distributed over a large number of specialized data bases (referenced in Table 1) dedicated to individual ncRNA families. A specialized data base for plant-specific ncRNAs is the *Arabidopsis* Small RNA Project Data base (ASRP) (Gustafson et al., 2005).

The purpose of this contribution is two-fold. Firstly, we tried to compile an overview of the current (January 2005) knowledge on all the different levels of RNA

activity in the cell, with an emphasis of what is (or is not) known about the evolution of individual classes of RNAs. Secondly, we use the framework of the review-like material to put new results on individual ncRNAs into perspective. Together, a picture emerges that on the one hand supports the picture of RNA as an ancient player in the cell, likely deriving from an RNA world pre-dating the last common ancestor of all extant life (Jeffares et al., 1998), while on the other hand many ncRNA families are probably relatively young innovations or have expanded dramatically, as for instance miRNAs, in certain lineages.

## Detection of ncRNAs

Genome data bases nowadays offer a wealth of annotation about protein-coding genes and their putative functions. Annotation of ncRNA genes, however, is almost non-existent. The main reason for this is the lack of established and reliable methods to detect such ncRNA genes computationally in genomic sequences. Current approaches for ncRNA detection can be clearly separated into two classes: methods to detect new members of already known and well-characterized ncRNA families, and attempts to predict RNA genes *de novo* so that novel families of ncRNAs can also be found.

### Members of known families

Large, highly conserved ncRNAs, in particular rRNAs, can easily be found using *blast* (Altschul et al., 1990). Similarly, *blast* can be used to find orthologous ncRNAs in closely related species, e.g. Tanzer and Stadler (2004) and Weber (2005). In most cases, however, this approach is limited by the relatively fast evolution of most ncRNAs. Since RNA sequence often evolves much faster than structure, the sensitivity of search tools can be greatly improved by using both sequence and secondary structure information.

The simplest class of search tools uses regular or context-free grammars to describe RNA motifs that are explicitly known to the user. There is no possibility to adapt the model to variations of the instance, and it is also very difficult for a user to define production rules for complicated motifs with a large number of exceptions.

With probabilistic models, such as stochastic context-free grammars (SCFG), the user is able to assign probability distributions to production rules; noise in the dataset is handled easily because the model can adapt itself to variations. The main drawback of stochastic context-free grammars is that most of the available implementations demand large computational resources. Hybrid languages, like HyPaL (Gräf et al., 2001) or the language used in RNAMotif (Macke et al., 2001), connect pattern languages with user-defined approximative rules, which rank the results according to their distance to the motif. Their advantage lies in a faster processing compared to SCFG. Nevertheless, the definition of approximative rules

also requires explicit knowledge, at least to some extent. Table 2 summarizes the most commonly used approaches.

**PatSearch** (Pesole et al., 2000b), **RNAMotif** (Macke et al., 2001), and **Palingol** (Billoud et al., 1996) are tools which allow the user to specify a given motif with a particular description language and offer search approaches to identify instances of the motif in a set of sequences.

**Palingol** is a constraint programming language to describe arbitrary rules on primary and secondary structure. The user defines a series of boolean expressions which must be satisfied by a successful hit.

In **PatSearch**, a language similar to regular expressions is used to describe motifs. For patterns composed of a string, a weight matrix can be defined which enables ranking and searching for approximative hits.

**RNAMotif** combines a pattern language with an awk-like programming language that describes approximative user defined scores. Sequences which have been matched successfully are evaluated and ranked according to the scoring section.

**ERPIN** (Gautheret and Lambert, 2001) is an example of tools that do not need an explicit definition of a descriptor to search for homologs of a motif. From a sequence alignment annotated with helix regions it extracts frequencies of nucleotides in single strands and base pair frequencies in helices. Those frequencies are compared to expected base frequencies in the target data base by calculating log-odds ratios. The sum of log-odds ratios over all positions of a target sequence gives the final score.

**RNAProfile** (Pavesi et al., 2004) requires as input the number of hairpins of a motif to extract it from an unaligned set of sequences where some contain the same motif. All sequences are folded and only those subsequences forming minimum free energy structures with the specified number of hairpins are regarded during the search. In a greedy search approach the selected regions of the first two sequences are pairwise aligned, according to primary and secondary structure. For each alignment, a profile, composed of observed frequencies of unpaired and paired nucleotides at each position, is defined and the best scoring ones are kept. In the second step the best scoring pairwise profiles are aligned to the selected regions of the next sequence and again only the best updated profiles are kept, and so on. If all sequences of the input set are processed, the highest scoring profiles define the detected motif. A fitness value is assigned to each final hit assessing its statistical significance.

A number of large-scale surveys have been performed using one of the general purpose tools mentioned above. A non-exhaustive list includes a miRNA survey using **ERPIN** (Legendre et al., 2005), a search for U5 snRNA and RNase P using **RNAMotif** (Collins et al., 2004), and a survey of RNase P RNAs in bacterial genomes (Li and Altman, 2004).

**Fragrep** (Mosig et al., 2004) is a simple sequence-based tool that allows to specify a query of short sequence elements that are separated by poorly conserved regions of variable length. Local alignment algorithms such as **blast** are therefore ill-suited for the discovery of new homologs of such ncRNAs in genomic sequences. The **fragrep** tool instead implements an efficient algorithm for detecting pattern fragments that occur in a given order. For each pattern fragment a mismatch

**Table 2.** General purpose algorithms for RNA motif detection (Tools that detect a special class of RNA motifs are not listed here)

Program	Comparative or single organism	Description
<b>Approaches which search for instances of a motif</b>		
ERPIN	Gautheret and Lambert (2001)	Input is a sequence alignment with consensus structure. For each helix and single strand a log-odds-score profile is defined which describes the motif
PATSearch	Pesole et al. (2000b)	Motif is defined by a language inspired by regular expressions
fragrep	Mosig et al. (2004)	Detects patterns consisting of approximately matched gapless blocks with constrained inter-block distances
RNAmot	Gautheret et al. (1990)	Finds combinations of primary or secondary structure motifs in nucleic acid sequences.
Palingol	Billoud et al. (1996)	A constraint programming language particularly adapted for secondary structures.
RNAMotif	Macke et al. (2001)	Allows both sequence and structure patterns, including pseudo-knots
infernal	Eddy (2002)	Description of structural motifs in terms of helices and sequence patterns. Putative hits are ranked according to user-defined rules
Rsearch	Klein and Eddy (2003)	Toolkit for constructing covariance models and finding new members of a family.
		Input is a multiple alignment with structural annotation. With SCFGs a consensus model of the RNA structure shared by these sequences is defined
		Input is a single RNA sequence and its structural information. Research is a local alignment algorithm which considers structural and sequence constraints. A base pair and single nucleotide substitution matrix for RNAs (RIBOSUM) defines alignment scores

**FastR** Bafna and Zhang (2004) Single Like Rsearch a pairwise alignment algorithm that addresses structural and sequence conservation. Running time is highly decreased by preprocessing the target sequences. Only those targets sharing similar structural features with the query RNA are aligned

#### Approaches which search for motifs from scratch

SLASH	Gorodkin et al. (1997)	Comparative	Input are unaligned sequences. foldalign defines highest scoring local alignments of these sequences according to sequence and structure constraints. COVE creates a SCFG model from those local alignments and does data base searches
RNAProfile	Pavesi et al. (2004)	Comparative	Input is a set of unaligned sequences. Motif is defined by the number of single hairpins it may contain. Greedy heuristic to find sequences in the input set which share a common motif with defined number of hairpins
GPRM	Hu (2003)	Comparative	Genetic programming approach to find structural RNA motifs that discriminate a set of input sequences from a set of randomized sequences
HyPa and HyPaLib	Gräf et al. (2001)	Single	A search engine and pattern library for “hybrid patterns”, consisting of sequence and structure elements. The language also includes thermodynamic constraints. Currently, however, HyPaLib contains only some 60 patterns

**Table 3.** Survey of Y RNAs in completely sequenced genomes using fragrep

Genome	Hs	Mm	Rn	Gg	Xt	Tr	Tn	Dr
# Matches	148	6	8	4	4	3	3	2

Hs: *Homo sapiens*, Mm: *Mus musculus*, Rn: *Rattus norvegicus*, Gg: *Gallus gallus*, Xt: *Xenopus tropicalis*, Tr: *Takifugu rubripes*, Tn: *Tetraodon nigroviridis*, Dr: *Danio rerio*.

tolerance and bounds on the length of the intervening sequences can be specified separately.

The application of fragrep is demonstrated in Table 3 using Y RNAs, an abundant small ncRNA described in some more detail below, as an example. It is straightforward to extract a query from sequences and structures of Y1, Y3, Y4 and Y5 RNAs given in O'Brien et al. (1993); the conserved sequence fragments of Y RNAs have also been studied by other authors (Farris et al., 1999; Teunissen et al., 2000). The large number of human sequences indicates that Y RNAs are associated with a repeat family in the human genome. An analysis of the Y RNA candidate sequences will be given in Section "Other snRNA-like molecules".

Specialized programs have been developed to detect members of particular ncRNA families. Examples of this approach include miRseeker for miRNAs (Lai et al., 2003), BRUCE for tmRNAs (Laslett et al., 2002), tRNAscan for tRNAs (Lowe and Eddy, 1997), snoScan for box C/D snoRNAs (Lowe and Eddy, 1999), fisher for box H/ACA snoRNAs (Edvardsson et al., 2003), as well as a heuristic for SRP RNAs (Regalia et al., 2002; Rosenblad et al., 2004). An improved method for box C/D snoRNAs was recently presented by Accardo et al. (2004): starting from yeast rRNA methylation sites, they first identified homologous positions in *D. melanogaster* rRNAs and then use snoScan (Lowe and Eddy, 1999) to search for putative snoRNAs with binding motifs complementary to the putative methylation sites. miRNAs in plants can be found by extracting those hairpin structures that contain sequence motifs complementary to a mRNA, which is then a putative target (Jones-Roades and Bartel, 2004; Bonnet et al., 2004a; Adai et al., 2005).

### Novel ncRNAs and RNA motifs

Detecting novel ncRNAs without any prior knowledge of sequence or structure is still a largely unsolved issue. In contrast to protein-coding genes, which show strong statistical signals like open reading frames or codon bias, ncRNAs lack any comparable signals in primary sequence that could be used for reliable detection.

Only in very special cases can ncRNAs be identified based on a significant bias in base composition. AT-rich hyper-thermophiles were successfully screened for ncRNAs simply by searching for GC-rich regions (Klein et al., 2002; Schattner, 2002). miRNAs can be detected based on their increased thermodynamic stability (Bonnet et al., 2004b). Carter et al. (Carter et al., 2001) used machine learning

techniques to extract common sequence features of known ncRNAs including GC content in *E. coli*.

Most ncRNAs do, however, depend on a well-defined structure for their function. This has led to various attempts to predict functional RNAs using predicted secondary structures. It was first suggested by Maizel and co-workers that functional RNA elements should have a more stable secondary structure than expected by chance (Le et al., 1988; Chen et al., 1990). However, Rivas and Eddy had to conclude in an in-depth study on the subject that thermodynamic stability alone is generally not statistically significant enough for reliable ncRNA detection (Rivas and Eddy, 2000). Some other characteristic measures derived from secondary structure predictions have been proposed (Schultes et al., 1999; Le et al., 2002, 2003) which, however, are also of limited value in the context of genome wide ncRNA prediction. A combination of gene expression data and high-level sequence conservation was successful in discovering novel ncRNAs in the intergenic regions of the *E. coli* genome (Wasserman et al., 2001).

The reason for the limited success of these approaches is that the presence of secondary structure in itself does not indicate any functional significance, because almost all RNA molecules form secondary structures. In fact, most compelling evidence for functional significance comes from comparative studies that demonstrate evolutionary conservation of structure.

Extensive computer simulations, see e.g. Schuster et al. (1994), Grüner et al. (1996a, b) and Huynen et al. (1996), showed that a small number of point mutations is very likely to cause large changes in the secondary structures. It follows that structural features will be preserved in RNA molecules with less than some 80% of sequence identity only if these features are under stabilizing selection, i.e., when they are functional.

This fact is exploited by the alidot (Hofacker et al., 1998) algorithm for searching conserved secondary structure patterns in large RNAs. Secondary structures are predicted independently for each sequence, typically using McCaskill's algorithm (McCaskill, 1990), which yields a list of thermodynamically plausible base pairs with their equilibrium probabilities. Next, a conventional multiple sequence alignment is computed, e.g., using ClustalW. By copying the gaps from the multiple sequence alignment into the predicted structures, a list of homologous base pairs is obtained. This list is then sorted by means of hierarchical credibility criteria that explicitly take into account both thermodynamic information and sequence covariation. A detailed description of the method can be found in (Hofacker et al., 1998; Hofacker and Stadler, 1999). A similar approach is taken by the ConStruct tool (Lück et al., 1996, 1999), which also features a graphical tool for manipulating the sequence alignment in order to achieve a better consensus structure. Alidot does not pre-suppose the existence of a global conserved structure. It is therefore particularly well suited when the sequences are expected to contain only small structurally conserved regions, as is the case for example in RNA viruses.

For predicting globally conserved structures a different technique, “folding the alignment”, may be preferred. Here, the folding algorithm itself is modified to work on a sequence profile, or multiple sequence alignment, instead of a single sequence.

The two best known implementations of this approach are `pfold` (Knudsen and Hein, 1999, 2003), and `RNAalifold` (Hofacker et al., 2002). `pfold` is based on an stochastic context-free grammar, and thus uses parameters derived from a training set. It also makes explicit use of a predicted phylogenetic tree. `RNAalifold`, on the other hand, uses the standard energy model for RNA secondary structures, augmented with a covariation term that rewards consistent and compensatory mutations. Thus, for identical sequences, it gives the same result as the single sequence prediction from `RNAfold`. With a few (or even just two) related sequences these programs achieve prediction accuracies much higher than prediction methods for single sequences. The approach is limited by the accuracy of the input alignment.

For sequences with less than 60% identity, pure sequence alignments typically differ significantly from structurally correct alignments. In these cases, one can resort to using a variant of the Sankoff algorithm (Sankoff, 1985) which computes the alignment and consensus structure simultaneously. Notable implementations are `foldalign` (Gorodkin et al., 1997, 2001b; Havgaard et al., 2005), `dynalign` (Mathews and Turner, 2002), `pmcomp/pmmulti` (Hofacker et al., 2004a), and `dart` (Holmes, 2004). The Sankoff algorithm is computationally very expensive, scaling as  $\mathcal{O}(n^6)$  in the unrestricted case. The above algorithms therefore use various restrictions to improve speed (`foldalign` for example considers only unbranched stem-loop structures). Nevertheless, they are generally not suitable for genome wide scans. A different approach to structural alignments is provided by making use of the tree representations of RNA secondary structures. Both `RNAforrester` (Höchsmann et al., 2003) and `MARNA` (Siebert and Backofen, 2003.) produce multiple alignments from pairwise structure-based alignments. For a recent comparison of techniques for consensus structure prediction see Gardner and Giegerich (2004).

Accurate predictions of consensus structures can provide a stepping stone towards reliable detection of functional RNAs. However, a successful ncRNA finder must also provide a measure of significance, such as an *p*-value or *E*-value. A well-known program to classify pairwise sequence alignments as ncRNA, protein coding, or anything else, is `qrna` (Rivas and Eddy, 2001). This program compares the score of three distinct models of sequence evolution to decide which one describes best the given alignment: a pair stochastic context-free grammars (SCFG) is used to model the evolution of secondary structure, a pair hidden Markov model (HMM) describes the evolution of protein coding sequence, and a different pair HMM implements the null model of a non-coding sequence. `Qrna` was successfully used to predict ncRNAs candidates in *E. coli* and *Saccharomyces cerevisiae* (Rivas et al., 2001; McCutcheon and Eddy, 2003), some of which could be verified experimentally. `Qrna` is, however, currently limited to pairwise alignments, and somewhat slow for large genomic scans. Other recent programs for detecting conserved RNA secondary structures include `ddbRNA` (di Bernardo et al., 2003) and `MSARi` (Coventry et al., 2004). A phylogenetic shadowing approach specifically geared towards the detection of miRNA precursors is described in Berezikov et al. (2005).

Currently, the sensitivity and/or specificity of all these programs is insufficient for screens of large eukaryotic genomes. Part of the problem is often oversimplification of the folding model (poor thermodynamics), as well as considering only

compensatory mutations as signal for structural conservation. Typical data sets, however, do not always show enough sequence variation ensuring this to be a significant indicator.

Recently, it has been demonstrated that the comparative approach can give significant results even for alignments with only few sequence and high similarity (Washietl and Hofacker, 2004). This approach uses RNAalifold (Hofacker et al., 2002) to compute consensus structures, making best use of covariance information and thermodynamic stability. Significance is then measured by a *z*-score comparing the consensus folding energy of the native alignment (as computed by RNAalifold) with the folding energies of randomized alignments, obtained by a shuffling procedure. Although the results are promising in terms of accuracy, the practicability of this approach is limited by the high computational costs caused by the time consuming shuffling procedure. In a more recent contribution, this problem is solved resulting in a time efficient algorithm showing similar accuracy. The program RNAZ (Washietl et al., 2005) uses two-independent criteria for classification: a *z*-score measuring thermodynamic stability of individual sequences, and a *structure conservation index* obtained by comparing folding energies of the individual sequences with the predicted consensus folding. The two criteria are combined by a support vector machine that detects conserved and stable RNA secondary structures with high sensitivity and specificity. Thus, RNAZ seems to be the first program suitable for screening large eukaryotic genomes (Washietl et al., 2005; Dieterich et al., 2005).

GPRM (Hu, 2003) considers motif prediction as a supervised learning problem. Coregulated mRNA sequences are used as positive examples, while the same number of randomly generated sequences form a set of negative examples. A genetic programming approach is used to learn the motifs in the predicted structures that can discriminate the positive set from the random sequences. Optimal discriminators are therefore good candidates for functionally important structural motifs (Hu, 2002).

It should be pointed out, however, that not all ncRNAs can be tracked down by searching for conserved secondary structures. To mention only a few examples, the U4 and U6 spliceosomal RNAs are known to form extensive *inter*-molecular interactions rather than forming stable *intra*-molecular secondary structures and are therefore missed by this approach. Also, most of the C/D-class snoRNAs lack an easily detectable secondary structure. Thus, while reliable structural RNA gene finding programs have come into reach, a general RNA gene finder remains elusive.

## Sequence evolution of ncRNA families

### ncRNAs and phylogenetic inference

While, as we have seen in the previous section, sequence information alone is in general insufficient to detect ncRNAs, it can be used very well to elucidate the

evolutionary relationships of these genes, at least within a given family of ncRNAs or RNA motifs. Since most known ncRNAs have evolutionarily conserved structures, however, they are only approximately described by models assuming independent evolution of sequence positions. A more accurate treatment explicitly takes into account that sequence positions that form conserved base pairs are highly correlated. Corresponding models of sequence evolution are described, e.g., in Schöninger and von Haeseler (1999), Knudsen and Hein (1999), Savill et al. (2001) and Otsuka and Sugaya (2003). The phase package (Jow et al., 2002; Hudelot et al., 2003) implements such a model and is specifically designed to infer phylogenies from RNAs that have a conserved secondary structure.

These secondary structures, however, have rarely been used in molecular phylogenetics so far. An exception is the investigation into the history of RNase P and RNase MRP RNAs by David Penny and co-workers (Collins et al., 2000). This study uses RNA editing distances (Shapiro and Zhang, 1990) implemented in Vienna RNA package (Hofacker et al., 1994) to show that “RNA secondary structure is useful for evaluating evolutionary relatedness, even with sequences that cannot be aligned with confidence”. More recently, cladistic analyses based on RNA secondary structure (Caetano-Anollés, 2002a, b) have demonstrated this point convincingly, in particular at the level of deep phylogenies.

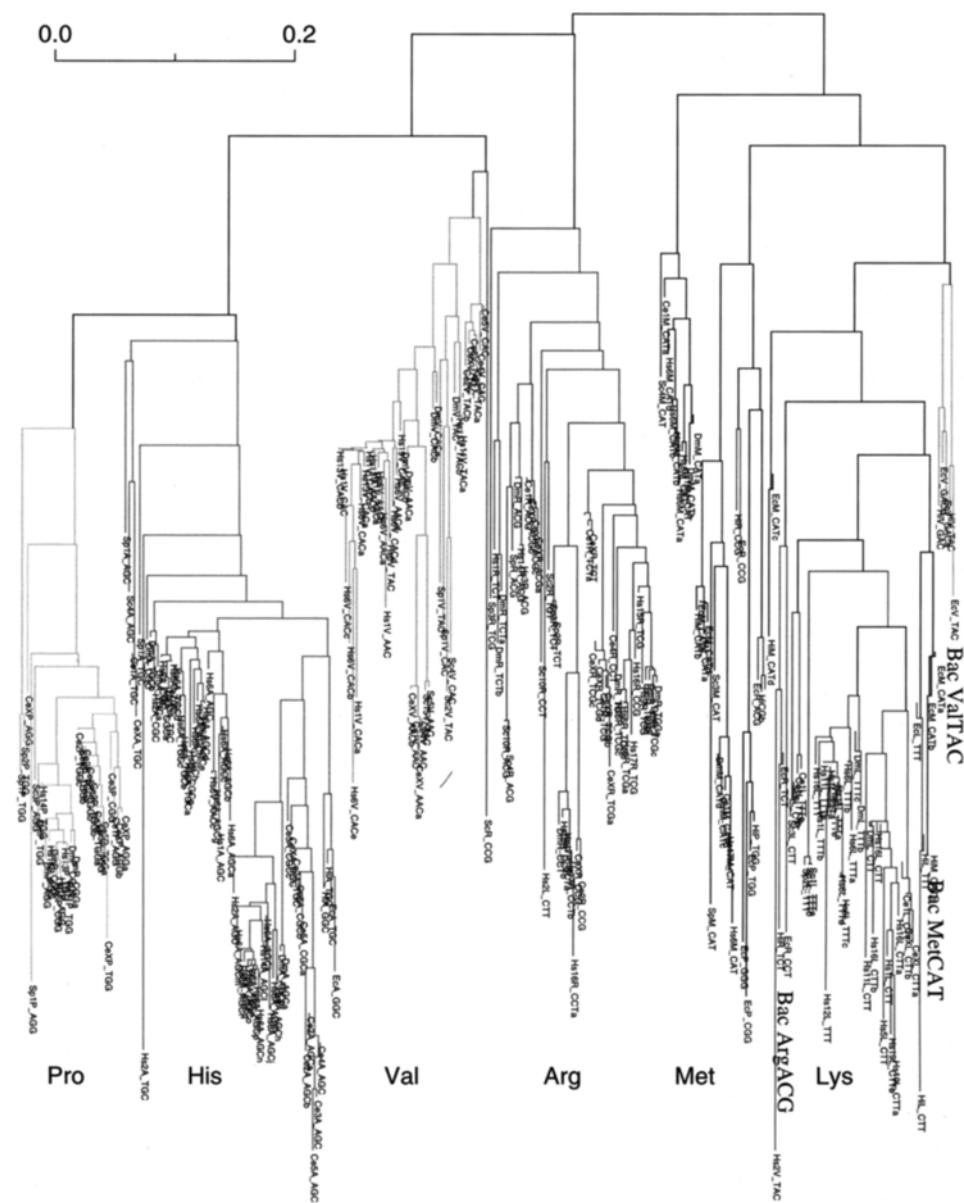
In the following, we compile an overview of our current knowledge of the evolution of the best known classes of ncRNAs. Our focus therefore are gene phylogenies and the history of duplications and losses that led to the present ncRNA inventory. This review of the literature is complemented by a number of original results, for which we provide supplemental data in electronic form.<sup>1</sup> We have mostly used neighbor-joining (Saitou and Nei, 1987) rather than the sophisticated maximum-likelihood techniques-mentioned above, since we are interested here in the large-scale patterns rather than subtle details of the ncRNA gene phylogenies.

## tRNAs

Multiple copies of functional tRNA genes, the existence of numerous pseudogenes and tRNA-derived repeats are general characteristics of tRNA evolution (Frenkel et al., 2004). Comparative sequence analysis of tRNA by means of statistical geometry provides strong evidence that tRNA sequences diverged long before the divergence of archaea and eubacteria (Eigen et al., 1989). In Fig. 1, we illustrate this using tRNAs coding for six of the 20 amino acids: tRNAs with the same anticodon form coherent subtrees. Models for the origin of tRNA from even simpler components are discussed, e.g., in Eigen and Winkler-Oswatitsch (1981), Rodin et al. (1993) and Di Giulio (2004).

The evolution of mitochondrial tRNA was studied in detail by Paul Higgs and collaborators (Hudelot et al., 2003; Higgs et al., 2003; Jameson et al., 2003). In particular, they present evidence that the two animal tRNA-Leu variants (one with anticodon UAG, the other with anticodon UAA) evolve by a peculiar mechanism of

<sup>1</sup><http://www.bioinf.uni-leipzig.de/Publications/SUPPLEMENTS/05-001/>



**Fig. 1.** Neighbor-joining tree (Saitou and Nei, 1987) of nuclear tRNAs with anticodons for Alanine (Ala), Valine (Val), Proline (Pro), Arginine (Arg), Leucine (Leu), and Methionine (Met) from Human, *Drosophila melanogaster*, *Caenorhabditis elegans*, *Schizosaccharomyces pombe*, *Saccharomyces cerevisiae*, *Escherichia coli* K-12, and *Haemophilus influenzae*. A few groups of bacterial tRNAs that fall outside the main groups are indicated. Sequence data are taken from the *Genomic tRNA Database*, <http://rna.wustl.edu/tRNADB/>.

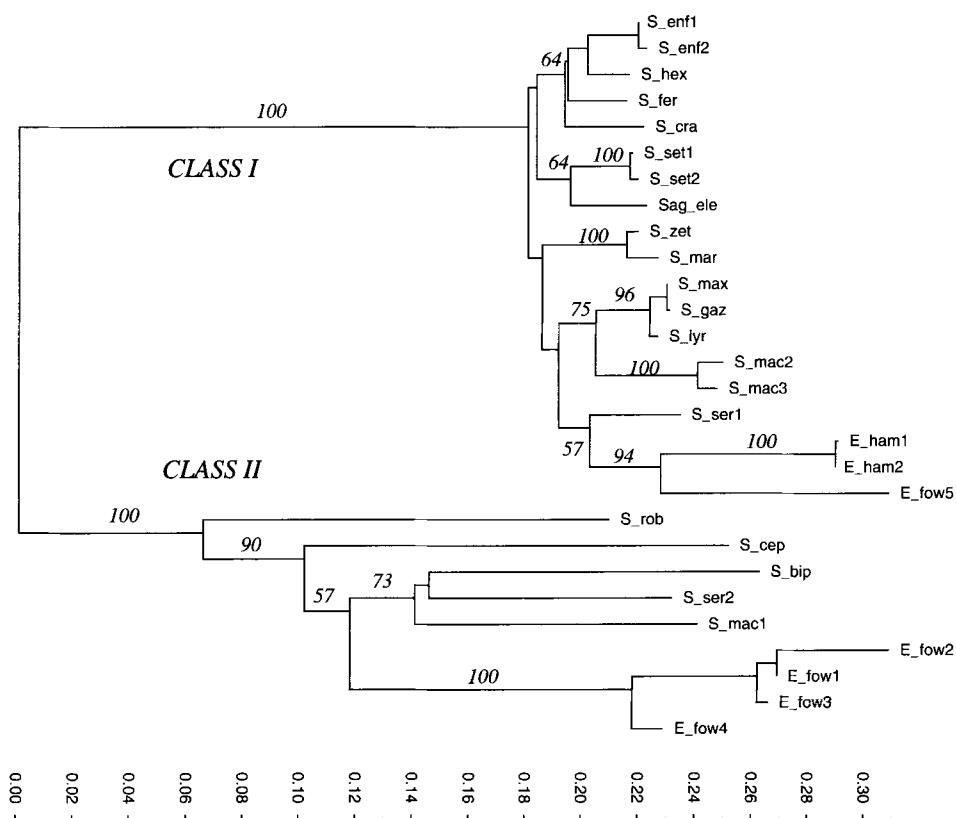
gene duplication, followed by mutation of the anticodon and subsequent gene loss. At least five such replacement events have been described in metazoan evolution (Jameson et al., 2003).

## Ribosomal RNAs

Evidence from both in vitro studies (Khaitovich et al., 1999; Nitta et al., 1998) and the analysis of the atomic structure (Ramakrishnan and Moore, 2001) reveals that the ribosome is in fact a ribozyme in which only rRNA is involved in the positioning of the A- and P-site substrates, and only RNA is in a position to chemically facilitate peptide-bond formation (Steitz and Moore, 2003). Due to its ubiquity, size, and generally slow rate of evolution, the small-subunit rRNA has become the most sequenced of all genes and an invaluable tool for molecular phylogenetics (Doolittle and Brown, 1994; Olsen and Woese, 1993; Van de Peer et al., 2000a; Peterson and Eernisse, 2001). More recently, large subunit rRNA are increasingly used for this purpose as well, e.g. Mallatt et al. (2004). The evolution of the secondary structures of rRNAs with an emphasis on functional sites is discussed in detail in Caetano-Anollés (2002b).

Most organisms have multiple copies of their rRNA genes. In *E. coli*, for instance, there are seven operons encoding rRNAs 16S, 23S, and 5S (Blattner et al., 1997). Typical Eukaryotes contain tandemly repeated arrays of rRNAs genes each of which contains three of the four rRNA components separated by two “internally transcribed spacers” (18S/ITS1/5.8S/ITS2/28S) (Hillis and Dixon, 1991). In most species the fourth rRNA gene, 5S rRNA, is also contained in this array, while it sometimes is dispersed throughout the genome (as in *Schizosaccharomyces pombe*, (Wood et al., 2002)), organized in its own tandem arrays (as in soybeans (Gottlob-McHugh et al., 1990)), or both (as in humans, Little and Braaten, 1989). For each of these genes, however, the rDNA sequences that are represented in fully processed rRNA are essentially identical in most organisms, i.e., rRNA genes are subject to *concerted evolution* (Hillis and Dixon, 1991; Schlötterer and Tautz, 1994; Gonzalez and Sylvester, 2001). This is the tendency of the different genes in a gene family or gene cluster to evolve “in concert”. As a consequence, one observes that paralogous sequence in the same species are more similar than orthologous sequences of different species. Multiple molecular mechanisms may account for this phenomenon: gene conversion (a non-reciprocal process in which two sequences interact in such a way that one is converted by the other), repeated unequal crossover, and gene amplification (frequent duplications and losses within family), see (Liao, 1999) for a review.

There are, however, exceptions to the rule: two classes of ancient paralogs of the 28S rRNA have been reported in the chaetognaths (Telford and Holland, 1997), see also Fig. 2. Similarly, paralog 18S rRNA are known, e.g., in the flatworm family *Dugesiidae* (Carranza et al., 1996, 1999) and in apicomplexans (Rooney, 2004), intraspecific 5.8S RNA variations have been reported in the coral *Acropora* (Márquez et al., 2003). In *Xenopus*, a somatic and an oocyte class of 5S RNA genes are differentially expressed in development due to changes in transcription factor



**Fig. 2.** Neighbor-joining phylogeny (Saitou and Nei, 1987) of partial 28S RNA sequences from chaetognatha. The tree is recalculated from data published by Telford and Holland (1997) using a clustalw alignment and the phylip package. The 28S sequences fall into two paralog groups that have separated at a common ancestor of the recent chaetognaths. For the species *Eukrohnia fowleri*, *Sagitta macrocephala* and *Sagitta serratodentata* both paralogs have been identified (Telford and Holland, 1997). Bootstrap values in percent (1000 replicates) are marked at major branches.

and histone interactions with the two types of gene (Wolffe, 1994). Distinct types of rRNA operons were also found in the *B. cereus* group (Candelon et al., 2004). Divergent paralogs could, if undetected, misguide phylogenetic studies.

### Spliceosomal RNAs

Most genes in higher eukaryotes contain introns that must be excised from the primary transcript to yield a mature mRNA. Intron removal and ligation of the exons occurs in a massive ribonucleoparticle (RNP), the *spliceosome*, see e.g. Nilsen (2003) and the references therein. Recently, there has been mounting evidence that main catalytic function in the spliceosome are indeed performed by its RNA

**Table 4.** Spliceosomal RNA components

Mechanism	snRNAs				
	Pol-II			Pol-III	
Major spliceosome	U1	U2	U4	U5	U6
Minor spliceosome	U11	U12	U4atac	U5	U6atac
Transsplicing		U2	U4	U5	U6

components, i.e., that the spliceosome, like the ribosome, is essentially a ribozyme (Valadkhan and Manley, 2001, 2003; Turner et al., 2004). The spliceosomal RNA U1 has an additional function in the regulation of transcriptional initiation (Kwek et al., 2002).

There are three distinct splicing mechanisms that are all dependent on a small set of RNA components of the spliceosome, Table 4: The major-spliceosome is the predominant mechanism, e.g., in vertebrates, plants, and yeasts, which spliced introns with the “canonical” GT–AG boundaries. The minor-spliceosome processes introns with non-canonical boundaries (Patel and Steitz, 2003), predominantly AT–AC. Trans-splicing, finally joins a small non-coding exon derived from the SL RNA to each coding exon of the pre-mRNA and is used to produce multiple mature mRNAs from a single poly-cistronic pre-mRNA (Pirotta, 2002; Tschudi and Ullu, 2002).

The evolutionary history of the spliceosome and its protein and RNA components is discussed in detail in Collins (2004). In spliced leader *trans* splicing, a common 5'-terminal exon is added post-transcriptionally to mRNAs which is derived from the SL RNA. The evolutionary origin(s) of this mechanism are still poorly understood because there is no clear pattern in the phylogenetic distribution of species that have this mechanism and the SL RNAs of distant species are too different to decide whether they are indeed homologous (Nilsen, 2001).

Both the pol-II transcribed spliceosomal RNAs U1, U2, U4, and U5 and the pol-III transcribed U6 snRNA appear in multiple copies in many vertebrates and are known to be subject to concerted evolution in some species (Domitrovich and Kunkel, 2003; Liao et al., 1997; Myslinski et al., 2004; Weiner and Denison, 1983). Divergent paralogs are also known in some species: for example, Xenopus has distinct embryonic and somatic classes of U1 snRNAs (Dahlberg and Lund, 1988). The evolution of U12 in vertebrates is considered in Tarn et al. (1995). A comprehensive investigation of snRNA evolution in the light of the available genomic sequence data, however, is still missing.

### Other snRNA-like molecules

**U7 RNA:** Replication-dependent histone pre-mRNAs, in contrast to all other mRNAs, are not polyadenylated. Instead, they are processed at their 3'end by endonucleolytic cleavage between two conserved sequence elements located within

**Table 5.** Repetitive elements associated with U7 snRNA

Species	Human	Mouse	Rat	Dog	Cow
# Hits	21/91 <sup>a</sup>	8	4	3	2

U7 RNA-like sequences are abundant in mammalian genomes, as determined by a *blast* search of the U7 sequence against the genomic sequence with a cutoff of  $E = 10^{-10}$ .

<sup>a</sup>Twenty-one hits when the U7 RNA sequence from Scharl and Steitz (1996) is used, 91 when using the consensus of all Rfam entries.

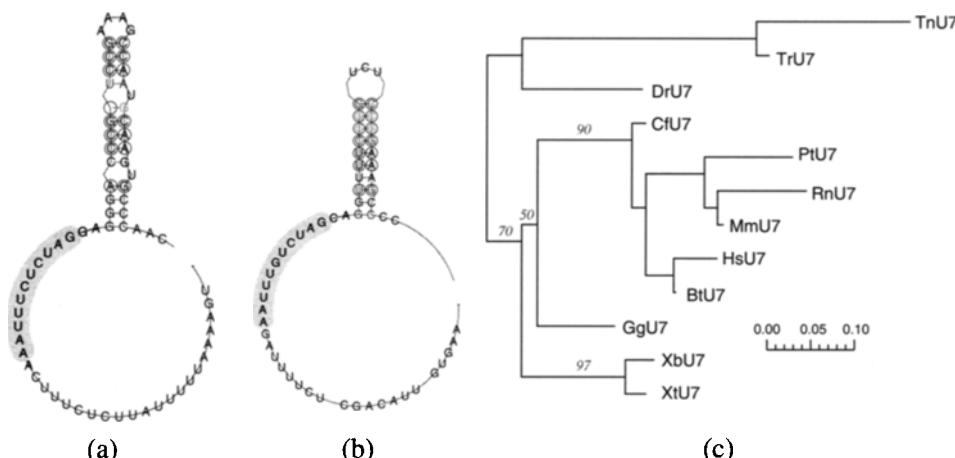
about 100 nt of the stop codon: a highly conserved stem-loop structure and a purine-rich histone downstream element (HDE). The latter is recognized by the U7 small nuclear ribonucleoprotein (snRNP) which consists of the U7 snRNA, a common Sm protein, and two unique Sm-like protein, as *Lsm10* and *Lsm11* (Schümperli and Pillai, 2004).

The U7-snRNP-dependent histone RNA 3'end processing mechanism is a metazoan innovation (Azzouz and Schümperli, 2003). Sequences of the U7 snRNA, which is only 60–70 nt long, have been published for some mammals (e.g. Soldati and Schümperli, 1988; Yu et al., 1996), Xenopus (Wu and Gall, 1993), Fugu (Myslinski et al., 2004), an echinoderm (Gilmartin et al., 1988), and more recently also for *D. melanogaster* (Dominski et al., 2003). Using a simple *blast* search, we found additional homologs in the chick genome, in two additional teleosts and in *D. pseudoobscura*. Like most other snRNAs, there are U7-derived repetitive sequences in some lineages, notably in human, while other species exhibit only a few scattered paralogs or pseudogenes (Phillips and Turner, 1991), or even have only a single copy (e.g. in the fugu (Myslinski et al., 2004)), see Table 5.

The sequences evolve quickly, severely limiting the power of comparative approaches. Because of the short sequence length of only 60–70 nt, one cannot expect a strong phylogenetic signal. Fig. 3c shows, however, that the sequence evolution is at least consistent with established phylogeny.

The U7 snRNA forms a relatively well-conserved hairpin structure just downstream of the Sm-binding sequence, see Fig. 3a and b. The U7 sequences were indeed used as an example to demonstrate the ConStruct approach to determining evolutionarily conserved secondary structures in Lück et al. (1999). The analysis in Fig. 3 using RNAAlifold (Hofacker et al., 2002) shows that there are significant differences in the secondary structures of invertebrates and vertebrates: vertebrate have smaller stem-loop structure with smaller or no interior loops or bulges.

**SRP RNA:** The signal recognition particle (SRP) is responsible for targeting nascent proteins to the ER membrane. In the process, protein synthesis is arrested when the SRP binds to the N-terminal signal of the nascent protein chain (Keenan et al., 2001). The SRP, components of which have been identified in all three domains of life (Rosenblad et al., 2003), contains a ncRNA, which in higher metazoan is also known as 7SL RNA. While the secondary structure of archaeal SRP RNAs closely resembles those of higher eukaryotes, Fig. 4, protozoan and fungal sequences may

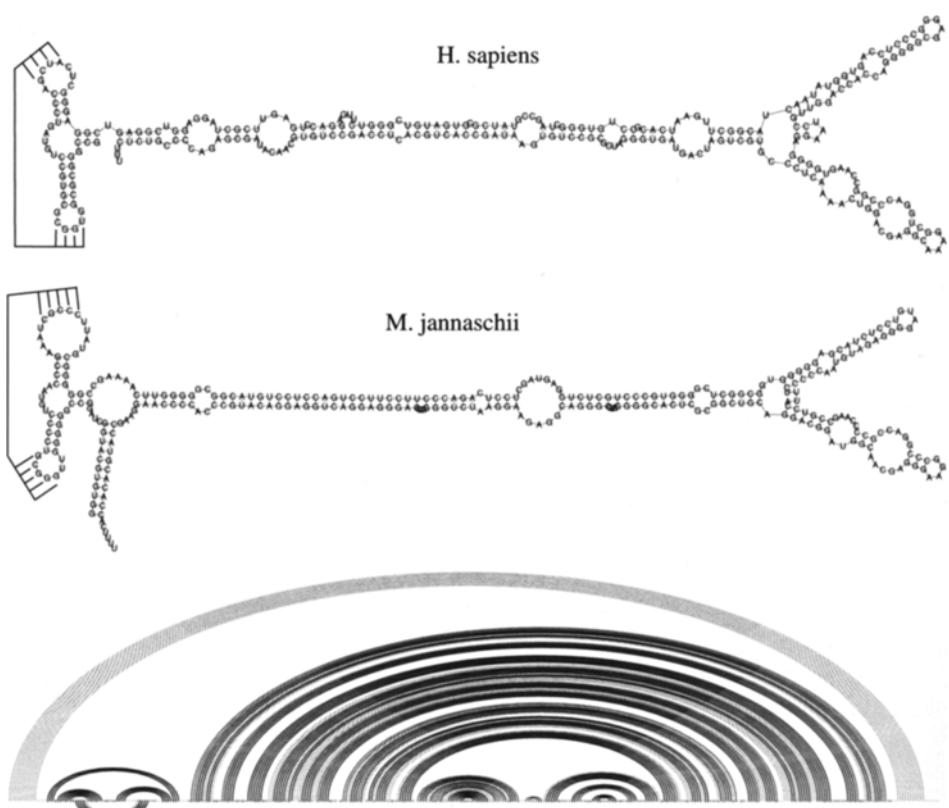


**Fig. 3.** Consensus secondary structures obtained from manual alignments of (a) 4 invertebrate and (b) 12 vertebrate U7 sequences calculated by RNAalifold (Hofacker et al., 2002). The highly conserved Sm-binding sequence is highlighted. Panel (c) shows a neighbor-joining tree obtained from the vertebrate alignment using the phylip package. Resolution within the mammals is poor, otherwise the U7 RNA tree reflects the accepted species phylogeny. Species abbreviations are: Bt *Bos taurus*, Cf *Canis familiaris*, Dr *Danio rerio*, Gg *Gallus gallus*, Hs *Homo sapiens*, Mm *Mus musculus*, Pt *Pan troglodytes*, Rn *Rattus norvegicus*, Tr *Tetraodon nigroviridis*, Tt *Takifugu rubripes*, Xb *Xenopus borealis*, Xt *Xenopus tropicalis*.

deviate considerably, and only the S-domain is present in most bacterial sequences (Zwieb and Eichler, 2002; Rosenblad et al., 2003, 2004). Chloroplast SRP RNA is described in Rosenblad and Samuelsson (2004). A detailed comparative discussion of the structural features of SRP RNAs from the different kingdoms can be found in (Zwieb et al., 2005).

Two small RNAs designated sRNA-85 (in *Leptomonas collosoma*, Ben-Shlomo et al. (1999)) and sRNA-76 (in *Trypanosoma brucei*, see Beja et al. (1993)) co-isolate with the 7SL RNAs of these Trypanosomatids, and there are indications that they function in place of certain protein components of the signal recognition particle. Their evolutionary relationship with the 7SL RNAs, however, is unclear (Zwieb et al., 2005).

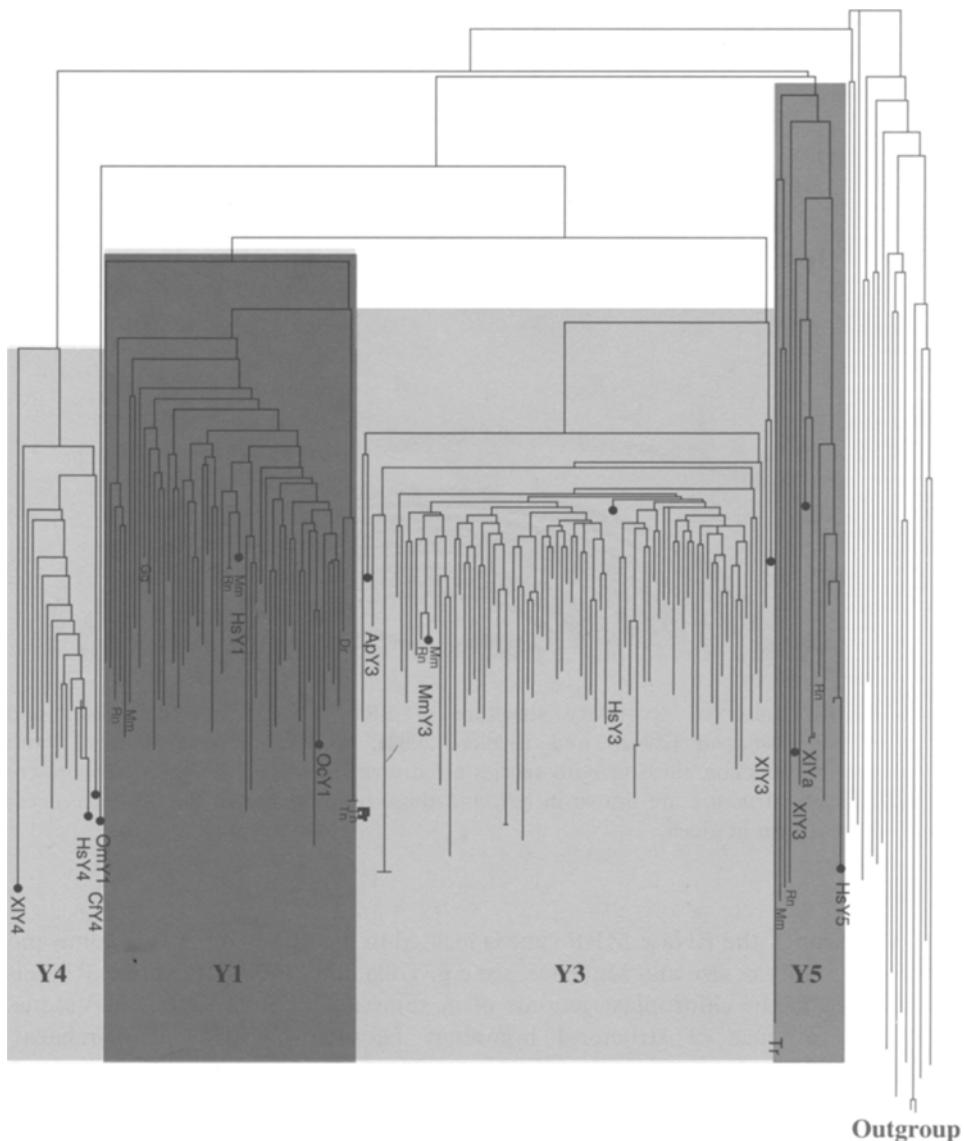
**P and MRP RNA:** The RNase P and RNase MRP RNAs are the catalytically active components of their respective RNPs, which both act as endonucleases. RNase P is essential for the maturation of tRNAs in Bacteria, Eukarya, and Archaea, see Pitulle et al. (1998) for a summary of its phylogenetic distribution and structural evolution. MRP RNA, in contrast, has been found only in Eukarya where it cleaves the primers necessary for the initiation of mitochondrial DNA replication (Morrissey and Tollervey, 1995), but also has nuclear functions. RNase P and MRP appear to be ancient paralogs, albeit it remains unclear whether MRP RNA is a eukaryote innovation or an older invention (Collins et al., 2000). In several



**Fig. 4.** Highly conserved secondary structure of SRP RNA from *H. sapiens* and *Methanococcus jannaschii* (Zwieb and Eichler, 2002). Bottom: superposition of both structures: base pairs contained in both species are drawn in black, base pairs only present in the *H. sapiens* structure are drawn in red, and those only present in the *Methanococcus jannaschii* are drawn in green.

ascomycete fungi, the RNase MRP gene is located in the mitochondrial genome and vary considerably in size and sequence, see e.g. Talla et al. (2005). RNase P RNA is also encoded in the chloroplast genome of some algae (de la Cruz and Vioque, 2003). The absence of structural homology between bacterial and archaeal/eukaryotic RNase P proteins suggests that RNase P once was a pure ribozyme that pursued completely different strategies in the recruitment of protein subunits in the two different lineages (Hartmann and Hartmann, 2003). A detailed investigation of bacterial RNase P RNAs (Haas et al., 1996) demonstrates an abrupt, dramatic restructuring in the common ancestor of the *Bacillus–Lactobacillus–Streptococcus* and the *Mycoplasma* groups of the low G+C Gram-positive bacteria. The latter shares the common ancestral “type A” structural architecture of bacterial RNase P RNAs, see also Krasilnikov et al. (2004) and Westhof and Massire (2004).

Expressed paralogs of RNase P RNA have been found in the mouse (Li and Williams, 1995), a systematic study of RNase P and MRP RNA variants, however, has not been performed to our knowledge.



**Fig. 5.** Neighbor-joining tree derived from the candidate Y RNA matches obtained by fragrep using a Clustalw alignment (Thompson et al., 1994). Known Y1, Y3, Y4, and Y5 candidate sequences were added to the candidate match sequences and are highlighted in the tree. Beside the outgroup on the left-hand side, all matching sequences can be clearly assigned to one of the known groups of Y RNA.

**7SK RNA:** Despite its abundance in mammalian cells, the function of the 7SK RNP has remained unknown until recent studies implicated 7SK RNA as well as components of the splicing apparatus (Kwek et al., 2002) in the regulation of transcriptional elongation, see Blencowe (2002), Michels et al. (2004) and Yik et al. (2003). Its secondary structure is known in detail from chemical probing experiments (Wasserman and Steitz, 1991). Interestingly, the 7SK RNA is very well conserved among vertebrates, while the lamprey sequence is already rather diverged (Gürsoy et al., 2000). D. Koper's Ph.D. dissertation (Koper-Emde, 2004) reports divergent 7SK sequences from the hagfish *Myxine glutinosa* and from two invertebrate species: *Branchiostoma lanceolatum* and *Helix pomatia*.

**Y RNAs** are small eukaryotic RNAs that are part of the Ro ribonucleoprotein (Ro RNP) complex, whose function is not known at present. Four families of Y RNAs, Y1, Y3, Y4, and Y5, have been described in human and frog. Their secondary structure is very well conserved among vertebrates (O'Brien et al., 1993; Farris et al., 1999; Teunissen et al., 2000). It consists of at least three stems, two of which form a stem-loop structure separated by a relatively short interior loop. The sequences in the stems, as well as parts of the loop regions, are highly conserved and probably serve as binding sites to the Ro60 protein in the Ro RNP complex and/or other cellular nucleic acids.

These conserved sequence patterns were used to scan genomic sequences for Y RNA candidates using the fragrep tool (Mosig et al., 2004), see Section "Members of known families". The phylogenetic tree resulting from an alignment of the matching sequences is shown in Fig. 5. It allows a further classification of the Y RNA candidate matches. Several matches, classified as an outgroup in the tree, are likely to be random occurrences of the search pattern. Integration of known representatives of the known classes of Y RNA (Y1, Y3, Y4, Y5) allows all other matches to be assigned to one of these known Y RNA classes. The data suggest that the four Y RNA families are at least as old as the last common ancestor of tetrapoda and actinopterygian fishes. The Y RNA family as whole is much older: a single member has been found in *Caenorhabditis elegans* (Van Horn et al., 1995).

**Vault RNAs** belong to a class of pol-III transcribed RNA genes with poorly understood function. Vaults are cytoplasmic ribonucleoprotein particles believed to be involved in multidrug resistance. The complex contains several small untranslated RNA molecules (van Zon et al., 2001). So far, vault RNAs have been described only for a few vertebrate species. Vault particles, however, are known also in the slime mold *Dictyostelium discoideum* (Vasu and Rome, 1995), suggesting that vault RNAs are at least as old as Eukaryotes. The human genome contains at least 4 distinct vaultRNA genes, three of which are located in small cluster and share external promoter elements (van Zon et al., 2001).

### **Small nucleolar RNAs (snoRNAs)**

Nascent rRNA transcript are matured in both eukarya and archaea (Dennis et al., 2001; Omer et al., 2000) with the help of a large number ribonucleoparticles that modify bases and direct cleavage. The human rRNAs, for instance, together contain

more than 200 modified nucleotides (Maden, 1990). The position of the snoRNA function is determined by the formation of a local snoRNA–rRNA duplex. Two major classes of snoRNA can be distinguished: the C/D box snoRNAs direct 2'-O-methylation of the ribose, while the H/ACA box snoRNAs guide the conversion of uridine nucleotides to pseudouridine. For details we refer to a series of reviews of snoRNA structure and function (Weinstein and Steitz, 1999; Kiss, 2001; Bachellerie et al., 2002; Terns and Terns, 2002; Henras et al., 2004).

Besides their canonical roles in rRNA maturation, snoRNAs also target spliceosomal RNA. These snoRNAs perform their function in the Cajal bodies; for this reason they are sometimes referred to as scaRNAs ("small Cajal-body associated RNAs") (Kiss, 2001). Most recently, three novel C/D box snoRNAs targeting U2, U4, and U12 snRNAs were identified that, in contrast to all other known metazoan snoRNAs, are independently transcribed (Tycowski et al., 2004). In archaea, tRNAs are also targeted for modification (Tang et al., 2002), in trypanosomatids the spliced leader SL RNA is modified as well (Liang et al., 2002; Uliel et al., 2004). An intriguing representative of this group is U85, a hybrid snoRNA that has both a functional C/D box and a functional H/ACA box domain that simultaneously modify the U5 snRNA (Jády and Kiss, 2001). Some snoRNAs lack complementarity to rRNAs or snRNAs. A small group of "orphan snoRNAs" (U3, U8, U22 and yeast snR10) directs rRNA cleavage instead of modification. The C/D box snoRNA U14, as well as the H/ACA box snoRNAs U17 (also called E1, and homologous to yeast sn30), E2 and E3, are both functional modification guides and play an additional role in pre-rRNA cleavage (Enright et al., 1996). An increasing number of recently identified snoRNAs exhibits tissue-specific expression patterns in contrast to all snoRNAs that are known to modify rRNA or snRNA (Cavaillé et al., 2000). The genes of these, mostly brain-specific, RNAs are subject to genomic imprinting. Vertebrate telomerase (Lue, 2004), finally, contains a conserved H/ACA box snoRNA domain (Mitchell et al., 1999; Chen et al., 2000).

The origin of the snoRNA machinery is still not well understood. The absence of snoRNAs from bacterial genomes suggests that snoRNPs arose in the archaeal and eukaryotic branch after the divergence of the bacteria. The K-turn motif, which forms the functional core of both classes of snoRNAs in archaea, on the other hand, also appears in bacterial RNAs including rRNA; it was probably present in the translation apparatus already before the last common ancestor (Penny and Poole, 1999). This suggests a common origin of both modern ribosome and modern snoRNPs from a primitive translation apparatus (Tran et al., 2004). The numerous box C/D and H/ACA RNPs of Archaea and Eukarya are likely to have arisen through duplication and variation of the guide sequence (Lafontaine and Tollervey, 2002). This scenario explains the lack of conservation of modified nucleotides shared between Archaea and Eukarya as well as the existence of tissue specific snoRNAs. In the following we demonstrate that this process is ongoing in vertebrate evolution.

The systematic investigation of snoRNA evolution is complicated by their fast evolution at sequence level. blast searches starting from human snoRNAs, for example, are usually unsuccessful already in non-mammalian vertebrate genomes. As a starting point for investigating the evolution of snoRNAs we have therefore

focused on the three snoRNAs that were first discovered (Nag et al., 1993), since sequences for these examples have been reported from a variety of different vertebrates. All three belong to the H/ACA class and are intron-encoded (Selvamurugan and Eliceiri, 1995; Mishra and Eliceiri, 1997).

The U17 (or E1) snoRNA is essential for the cleavage of pre-rRNA within the 5' external transcribed spacer (ETS) (Enright et al., 1996) with a length ranging from 200–230 nt, longer than most snoRNAs; its secondary structure has been studied in detail (Cervelli et al., 2002). Its sequence evolution in chelonians is discussed in (Cervelli et al., 2003). Both E2 and E3 snoRNA are involved in the processing of eukaryotic pre-rRNA and have regions of complementarity to 28S rRNA. Gene trees reconstructed for these three examples are displayed in Fig. 6. While in many cases closely related paralogs are found, we can also identify ancient duplications that have been maintained in the genome over long times.

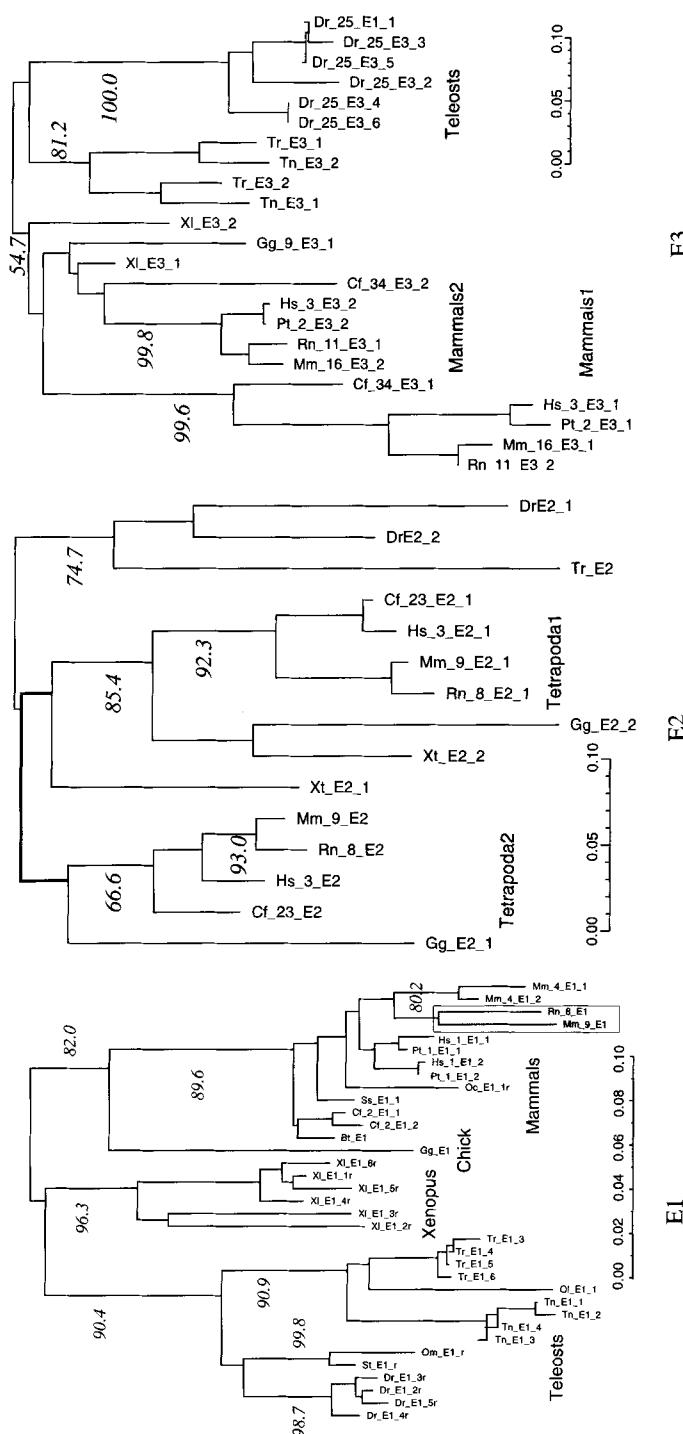
The evolutionary history of the paralog snoRNAs differs considerably between the three examples. The U17/E1 sequences for each species cluster together (with the exception of the Human and Chimp sequences), suggesting that the paralogs (which reside in adjacent introns) evolve via concerted evolution. In addition, however, the rodent genomes contain an additional paralog located on a different chromosome. In contrast, both for E2 and E3 we find two distinct evolutionary old paralog groups. In the case of E2 they separated before the advent of the tetrapods; the split between the two E3 groups predates the last common ancestor of the eutherian mammals. The six copies of E3 in the zebrafish apparently arose after the teleost-specific genome duplication (Amores et al., 1998).

The history of only a few other snoRNAs has been investigated in detail. Maybe the most interesting example is the C/D box snoRNA U36. It is homologous to snR47 in *S. cerevisiae* and appears in two paralogs in adjacent introns of the rpL7a gene in non-mammalian vertebrates. In mammals, however, U36a was duplicated with subsequent differential loss of function (Gilley and Fried, 1998). Other examples of snoRNAs whose evolution has been discussed in the literature include U14 (Samarsky et al., 1996) and U24 (Gilley and Fried, 1998).

The patterns observed in Fig. 6 show that concerted evolution breaks down occasionally when two paralogs acquire functional or regulatory differences. The mechanism behind the concerted evolution of snoRNA copies is not well known. The identification of a retrogene with a poly(A) tail for H/ACA box snoRNA U99 (Vitali et al., 2003), supports the idea that retro-transposition events play a substantial role in the mobility and diversification of snoRNA genes during evolution. This would argue for gene amplification (Weiner and Denison, 1983).

## Telomerase RNA

Telomeres are specialized protein-DNA complexes that cap chromosome ends that are essential for genome stability and cellular proliferation (Ferreira et al., 2004). Sequence loss during replication is counteracted by specialized mechanism(s) in organisms with linear chromosomes (Lingner et al., 1995). In most organisms, the



**Fig. 6.** Neighbor-joining trees of the E1, E2, and E3 snoRNAs. Bootstrap values from 1000 replicates are indicated in italics. The U17 sequences of *Takifugu rubripes* are taken from Acc. No. X94942 (Cecconi et al., 1996); Tr\_E1\_4 does not map unambiguously to a genomic location. The copies of the E1 snoRNAs that are located in a different host gene in rodents are highlighted.

telomerase RNP extends chromosome ends by iterative reverse transcription of its RNA template, the telomerase RNA (Kelleher et al., 2002).

The secondary structures of the telomerase RNAs from vertebrates, ciliates, and yeast vary dramatically in sequence composition and in their size but share a common core structure (Chen and Greider, 2004; Dandjinou et al., 2004; Lin et al., 2004; Zappulla and Cech, 2004) that hints at an ancient origin. Plants also contain well-conserved telomerase, see Oguchi et al. (2004) and the references therein; plant telomerase RNA, however, does not seem to have been studied systematically so far.

The vertebrate telomerase RNA apparently has co-opted a H/ACA box snoRNA domain (Mitchell et al., 1999) during its evolution, shares evolutionarily conserved proteins with H/ACA snoRNPs, and contains a Cajal body-specific localization signal that is shared with a Cajal body-specific subclass of H/ACA snoRNPs (Jády et al., 2004).

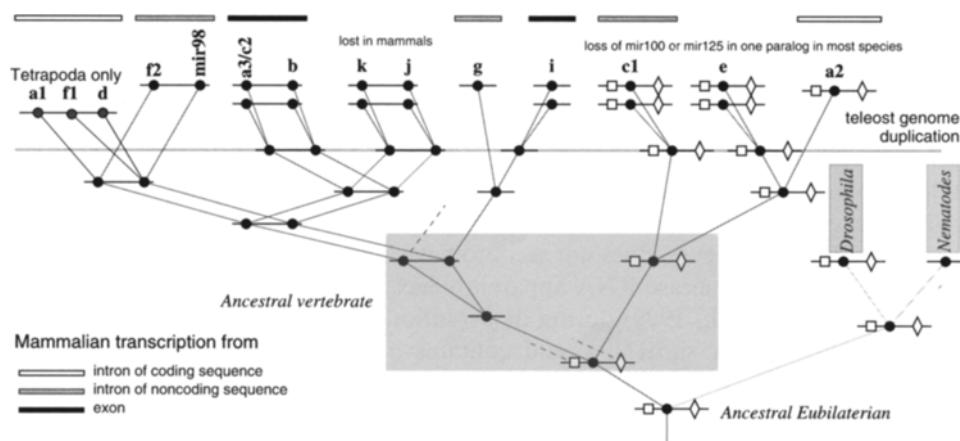
## MicroRNAs

miRNAs form a class of ncRNA genes whose products are small single-stranded RNAs with a length of about 22 nt. These are involved in the regulation of translation and degradation of mRNAs. We refer to the recent review (Nelson et al., 2003) for a discussion of their functions and mechanisms as well as their history of discovery. miRNAs are known in both multi-cellular animals and plants. A dedicated data base, the miRNA Registry (Griffiths-Jones, 2004), at present<sup>2</sup> contains more than 1345 miRNA sequences from 12 species. Recently, several miRNAs were detected using the micro-array technique (Barad et al., 2004; Sun et al., 2004).

Many of the known miRNAs appear in clusters on a single poly-cistronic transcript (Lee et al., 2002; Mourelatos et al., 2002; Lagos-Quintana et al., 2003; Lai et al., 2003). The **mir-17** family, for instance, consists of numerous paralogs of three apparently non-homologous sequences. A detailed investigation of its evolutionary history (Tanzer and Stadler, 2004) revealed a complicated sequence of tandem duplications within a cluster and duplications of entire clusters, which are probably linked to genome-wide duplications (Holland et al., 1994; Panopoulou et al., 2003). Two miRNA families that are associated with the *Hox*-clusters have received considerable attention: *mir-10* and *mir-196* (Yekta et al., 2004). Again, an expansion of both families is observed that closely follow the vertebrate and teleostean genome duplications (Tanzer et al., 2005).

As a further example we consider the history of the *let-7* family. The *let-7* gene was discovered in the *C. elegans* as a regulator in developmental timing (Reinhart et al., 2000). The *let-7* miRNA is present in diverse animal phyla including chordates, echinoderms, molluscs, annelids, arthropods, nematods, chaetognaths, nemerteans, and platyhelminths, but it is absent in basal metazoa including cnidarians, poriferans, ctenophora, and acoel flatworms (Pasquinelli et al., 2000, 2003). In

<sup>2</sup>Release 5.0, September 2004

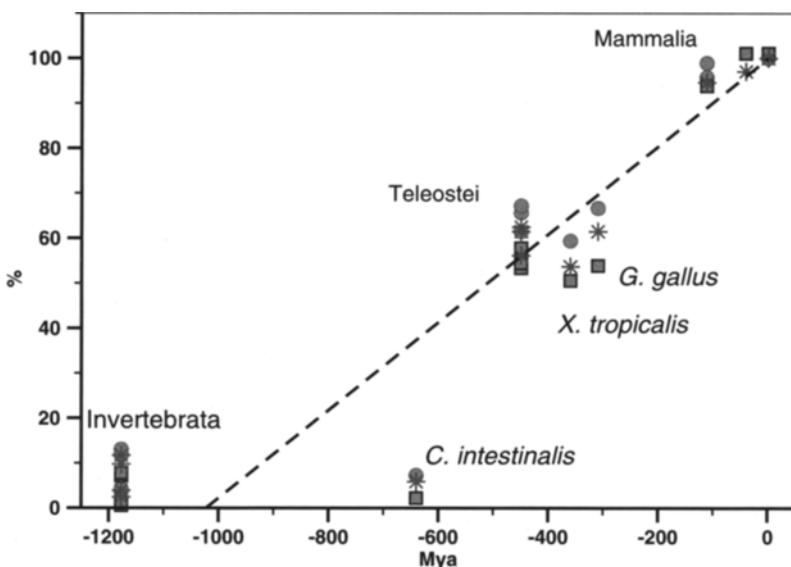


**Fig. 7.** A scenario for the expansion of the *let-7* family in vertebrates. There are few lineage-specific changes since the last common ancestor of teleosts and tetrapods; the inventory of *let-7* paralogs in tetrapods has essentially remained the same with the exception of the triple *let-7-a1*, *let-7-f1*, and *let-7-d1*. In teleost fishes we observe the loss of some paralog clusters and in particular the loss of linked *mir-100* ( $\square$ ) and/or *mir-125* ( $\diamond$ ) copies in some of the sequences of the *let-7-a2/c1/e* families. In addition, the *let-7-j/k* pair has been deleted in mammals (in contrast to birds). The bars in the top line indicate the mode of transcription of the human and mouse sequences (Rodriguez et al., 2004). A recent computational survey (Legendre et al., 2005) produced distant relatives of the vertebrate *let-7* sequences in *Ciona intestinalis* and *Ciona savignyi*. The sequences are too diverged, however, to determine whether they were produced by independent duplication of the ancestral *let-7* or whether they share part of the duplications shown in the figure.

vertebrates a plethora of *let-7* paralogs are known. In Fig. 7, we present a reconstruction of the history of this miRNA family.

Mammals seem to share a more or less similar miRNA repertoire. More than 90% of the mammalian miRNAs listed in the Rfam miRNA registry v4.0 (Griffiths-Jones et al., 2003) can be found in human, mouse, and rat. In contrast, chicken and frog contain only 50–60% of the mammalian miRNAs, whereas teleost fishes harbor slightly more (50–65%). Since the chicken and frog genome sequencing and assembly is still incomplete, these numbers might change slightly in future studies. If the number of miRNAs increased linearly in evolution, *Ciona intestinalis*, an ascidian urochordate and hence close relative of the vertebrates, would be expected to contain about 30% of the miRNAs found in mammals. However, we were able to detect only about 15%. This suggests that the origin of vertebrates was associated with a dramatic expansion of the miRNA repertoire.

Fig. 8 summarizes statistical evidence. Only a small group of miRNAs, which includes *let-7* (discussed above), *mir-10* (Tanzer et al., 2005), and *mir-92* (Tanzer and Stadler, 2004), can be found throughout most metazoans. These three families are characterized by numerous paralog miRNA genes at dispersed genomic locations



**Fig. 8.** Non-linear increase of miRNAs in Evolution. The human (green circles), rat (blue squares) and mouse (red stars) miRNAs listed in the Rfam mirNA registry v4.0 were blasted (cut off 10e-4) against the genomes of Invertebrata (*D. melanogaster*, *A. gambiae*, *C. elegans*, *C. briggsae*), *Ciona intestinalis*, Teleost fishes (*D. rerio*, *T. rubripes*, *T. nigroviridis*), *X. tropicalis*, *G. gallus* and Mammalia (*H. sapiens*, *M. musculus*, *R. norvegicus*). The percentage of mammalian miRNAs recovered are plotted against the evolutionary distance of those species.

and an additional expansion of families in teleosts. This points at a close association of the miRNA expansion with the genome duplications at the root of vertebrate tree (Lynch and Conery, 2000; Escriva et al., 2002) and early in the evolution of the actinopterygian lineage (Amores et al., 1998).

The origin of miRNAs remains unknown. As yet, no miRNA with homologs in both animals and plants has been described so far, although the miRNA processing machinery is clearly homologous. In Tanzer and Stadler (2004), it has been argued that miRNA could easily arise *de novo* since stem-loop structures resembling pre-miRNAs are very abundant secondary structures in genomic sequences. Most recently, a mechanism for the origin of new miRNAs in plants from inverted duplications of expressed sequences has been proposed for the *A. thaliana* sequences *mir161* and *mir163* (Allen et al., 2004). In this scenario, the new miRNAs will target the mRNA they arose from. On the other hand, evolutionarily ancient miRNAs are also known in plants: miR166 is conserved between flowering plants, ferns, mosses, and hornworts. In addition to land-plants and metazoan animals, miRNAs have also been found in viral genomes, including the Epstein–Barr virus (Herpesviridae) (Pfeffer et al., 2004) and HIV (Retroviridae) (Bennasser et al., 2004; Omoto et al., 2004).

## Other classes of small ncRNA

RNA editing in trypanosome mitochondria is a unique post-transcriptional maturation process in which uridine residues are inserted and/or deleted at precise sites of mitochondrial mRNAs (Brennicke et al., 1999; Estévez and Simpson, 1999; Gott and Emeson, 2000; Stuart et al., 1997). Guide RNAs (*gRNAs*), which are usually transcribed from the kinetoplast DNA minicircles (Hong and Simpson, 2003), provide the information for the editing.

In contrast, RNAediting mechanism (besides those the snoRNA-based base modifications) in other eukaryotes, prokaryotes, and viruses do not make use of RNA components (Korencic et al., 2004; Bishop et al., 2004). Models for the evolution of the gRNA-based editing process are discussed in Landweber (1992), a phylogenetic analysis of U-insertion editing (Landweber and Gilbert, 1994) suggests that extensive editing is a primitive genetic phenomenon that has disappeared in more modern organism, see also Simpson et al. (2000).

Probably the best-understood bacteria-specific ncRNA is the *tmRNA*, which is part of a ribonucleoprotein complex and combines the functions of tRNAs and mRNAs in order to rescue stalled ribosomes (Haebel et al., 2004). Usually tmRNA is a single molecule. At least three isolated clades in alpha-proteobacteria (Keiler et al., 2000), cyanobacteria (Gaudin et al., 2002; Williams, 2002), and beta-proteobacteria (Sharkady and Williams, 2004) have two-component tmRNAs, while jakobids have lost the mRNA-like region in their mitochondrial tmRNAs (Jacob et al., 2004). Reduction of the tmRNA structure in endosymbionts seems to be a common phenomenon (Gueneau de Novoa and Williams, 2004). The usefulness of tmRNA sequences for bacterial phylogenetics is demonstrated in Felden et al. (2001) by revealing a structural feature that is characteristic for beta-proteobacteria.

Prokaryotes contain a diverse set of small non-coding *sRNAs*. For example, a number of small (40–400 nt) RNAs that neither encode proteins nor function as tRNAs or rRNAs, have been characterized in *E. coli* (Hershberg et al., 2003; Vogel et al., 2003). The functions of many of these RNAs remain to be determined, while some of them are known to play crucial regulatory roles. There appear to be three general mechanisms: some are integral parts of RNP complexes, such as the 4.5S component of the signal recognition particle and RNase P RNA. A few, such as the 6S RNA, which regulated RNA polymerase activity (Montzka Wassarman and Storz, 2000), and the CsrB and CsrC RNAs mimic the structures of other nucleic acids, while a third class, reviewed in Storz et al. (2004), acts by specific base-pairing with other RNAs. The co-evolution of the small RNA *micF* and its target mRNA *ompF* in Enterobacteria was studied in some detail (Delihas, 2003). A curious case are the MCS4 RNAs in mycoplasmas, which have a sequence similarity with eukaryotic U6 snRNAs. Homologs in other bacteria do not seem to exist (Ushida et al., 2003), so that horizontal gene transfer from the host organism is a plausible explanation. Otherwise, very little is known about the origin and evolutionary relationships of the small ncRNAs in prokaryotes (Gottesman, 2004).

An increasing number of *viral* ncRNAs have been reported as well. Examples include the recently discovered viral miRNAs (Bennasser et al., 2004; Omoto et al.,

2004; Pfeffer et al., 2004), the well-known VA1 RNA of adenoviruses (Mathews, 1995), which is capable of inhibiting RNAi in human cells (Lu and Cullen, 2004), the pRNA component of the packaging motor in some bacteriophages (Bailey et al., 1990; Guo, 2002). One might suspect that at least some of the conserved RNA structure elements that were discovered in computations surveys of RNA virus genomes (Hofacker et al., 2004b; Thurner et al., 2004; Witwer et al., 2001) are also ncRNAs rather than *cis*-acting elements.

### mRNA-like ncRNAs

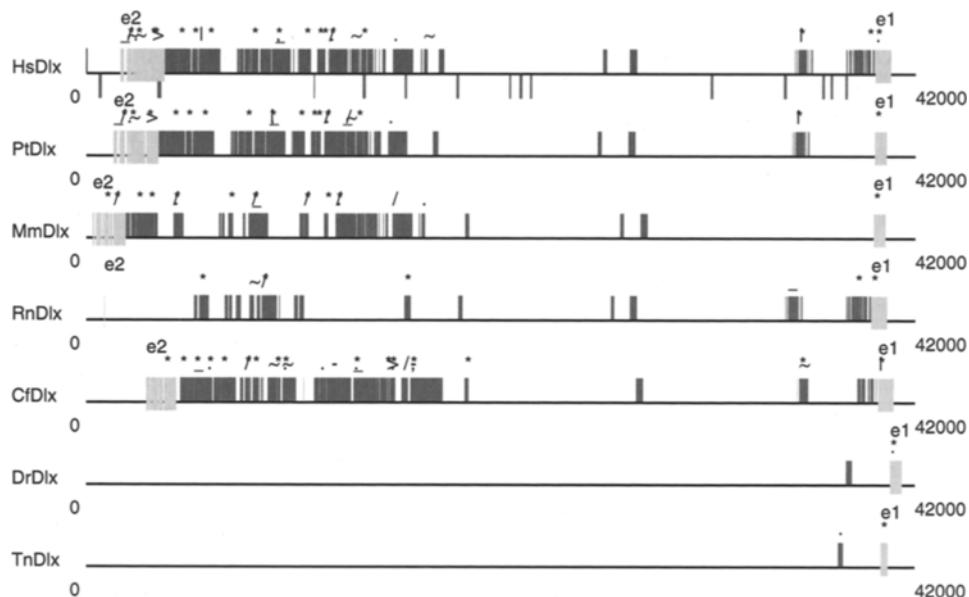
In eukaryotic cells, many RNA transcripts can be found that are not translated into protein. These so-called mRNA-like RNA transcripts are polyadenylated and spliced. In contrast to translated genes, they lack long ORFs (Erdmann et al., 1999, 2000). The best-known mammalian representatives of this rapidly expanding group are *H19* and *Xist*. Some of these large ncRNAs, including mammalian *Xist* and *Air*, and *roX* in *Drosophila*, have distinct roles in epigenetic gene regulation they are performed by means of chromatin modifications, reviewed in Andersen and Panning (2003). A number of plant specific mRNA-like ncRNAs are known experimentally; additional candidates were detected in a computational survey<sup>3</sup> of *A. thaliana* ESTs (MacIntosh et al., 2001).

The *Xist* (X-inactive specific transcript) gene is the only gene known to be specifically transcribed from the inactive X chromosome in female somatic cells (Brown et al., 1991). It codes for a 17-kb spliced, polyadenylated ncRNA. *Xist* is necessary and sufficient for the initiation and spread of X inactivation (Penny et al., 1996). The *Xist* gene is associated with an anti-sense transcript *Tsix* (Lee et al., 1999; Mise et al., 1999) that is thought to be a repressor for *Xist*. A comparative analysis of the X-Inactivation Center (XIC) region, and the *Xist* gene in particular, in human, mouse, and cow is reported in Chureau et al. (2002): while the *Xist* gene is well conserved among mammals with minor difference in the intron-exon structure, there is no apparent sequence conservation for the anti-sense transcript *Tsix*. Chureau et al. (2002) also identified two new ncRNA gene, termed *Jpx* and *Ftx* in the XIC region, which are well conserved in mammals.

The human *H19* gene is an imprinted gene that is exclusively expressed from the allele of maternal origin. It has a conserved secondary structure in mammals (Juan et al., 2000). The *H19* gene is abundantly expressed in both extraembryonic and fetal tissues and is repressed after birth, except in a few adult organs. The possible functional relationship between *H19* expression and tumorigenesis is still a matter of debate, as it seems to depend on the organ, the cell type and the cellular environment, see e.g. Berteaux et al. (2004) and the references therein.

As a third example, we describe here a computational analysis of the recently discovered mRNA-like ncRNA *evf-1*, which is located upstream of the *Dlx6* gene and its expression is linked to both Sonic hedgehog (*shh*) and *Dlx* genes (Kohitz and Fishell, 2004). *Dlx6* occurs clustered with *Dlx5*, another member of the same class of homeodomain transcription factors that are involved e.g. in the patterning and

<sup>3</sup><http://www.prl.msu.edu/PLANTncRNAs/>



**Fig. 9.** Sequence conservation in the region 42 kb upstream of the *Dlx6* gene. Boxes above the line mark phylogenetic footprints detected by *tracker* (Prohaska et al., 2004). Conserved regions that lie in the known exons of *evf-1* are colored green and cross the line. Boxes below the line for the human sequence mark rRNAZ hits. Putative *cis*-acting elements as identified using *infernal* and the Rfam data base are denoted by symbols above the *tracker* hits using the following symbols: \* IRE; \_ Hammerhead\_1; / SECIS; ~ REN-SRE; . Histone3; > U36; — Intron\_gpII; , s2m; - tRNA. The following sequences were used: HsDlx *Homo sapiens*, PtDlx *Pan troglodytes*, MmDlx *Mus musculus*, RnDlx *Rattus norvegicus*, CfDlx *Canis familiaris*, DrDlx *Danio rerio*, TnDlx *Tetraodon nigroviridis*.

migration of ventral forebrain neurons, see Sumiyama et al. (2003). Like *Xist* and *H19*, *evf-1* shows no homology to other known ncRNA sequences (Kohtz and Fishell, 2004).

The *evf-1* genes consists of two exons that are divided by a single approximately 37.5 kb large intron. We analyzed the DNA-sequence 42 kb upstream of the *Dlx6* gene to find highly conserved regions in this genomic region, Fig. 9. The highly conserved regions were detected by *tracker* (Prohaska et al., 2004), a program for phylogenetic footprinting (Zhang and Gerstein, 2003). The so detected phylogenetic footprints were scanned for conserved RNA-secondary structures using *RNAfold* and *alidot* (Hofacker et al., 1994, 1998; Hofacker and Stadler, 1999; Hofacker, 2003) and assigned to known secondary structure elements according Rfam (Griffiths-Jones et al., 2003) using S. Eddy's *infernal* program (Eddy, 2002). *Infernal* suggests a possible annotation for 34 of the 79 *tracker* hits. A table listing all blocks of conserved sequence elements can be found in the electronic supplement, see also Fig. 9. The position of the two exons was inferred by *blast* comparison with the rat sequence (Acc. no. AY518691.1).

In contrast to the mammalia-specific genes such as *Xist* and *H19* we find that *evf-1* shares at least *exon-1* and one large intronic sequence element with teleost fishes. A blast search also recovers *exon-1* from the *xenopus* and chicken genome. Since the genome assemblies of both the frog and the chick are incomplete in this region these sequences were not included in the analysis summarized in Fig. 9.

### Antisense RNAs

Antisense RNAs predominantly act as post-transcriptional downregulators of gene expression (Lavorgna et al., 2004). Indeed, some of the RNA families discussed above can be viewed as antisense RNAs since they exert their function by binding complementarily to their target RNAs; examples are the miRNAs, snoRNAs, as well as many of the bacterial small RNAs (Wagner and Flärdh, 2002). The analysis of genomic sequence data, however, has revealed that a substantial fraction of transcribed DNA does not code for proteins and often derives from the anti-sense strand, see e.g. Kampa et al. (2004), Shendure and Church (2002) and Yelin et al. (2003). Antisense transcripts thus emerge as a common mechanism of regulating gene expression in eukaryotic cells, reviewed e.g. in Lavorgna et al. (2004). Mechanistically, there are three major pathways: The formation of *double-stranded RNA* may trigger the RNAi pathway and lead to degradation of the sense transcript (Hannon, 2002). Binding of sense and anti-sense transcript may prevent the binding of other *trans*-acting factors (*RNA masking*). *Transcriptional interference* is the inhibition of transcriptional elongation due to a collision of the RNA Pol-II complexes on overlapping transcriptional units located at opposite strands (Precott and Proudfoot, 2002). Antisense RNAs are transcribed either *in cis* from the opposite strand, or *in trans* from a different genomic locus.

Many anti-sense transcripts are only poorly conserved in evolution, e.g. the *tsix* gene, which is the antisense transcript to the *Xist* ncRNA associated with X chromosome inactivation (Section “mRNA-like ncRNAs”). On the other hand, a number of well-conserved antisense transcripts are known. Probably the best-studied example is the *HoxA11* antisense transcript, which is well-conserved between human and mouse and exhibits tissue-specific alternative splicing (Potter and Branford, 1998). The Na/Pi co-transporter is essential in maintaining phosphate homeostasis in vertebrates. Antisense transcripts associated with the *npt* genes have been described in wide range of vertebrates (Werner et al., 2002) suggesting a conserved mode of transcription. Natural anti-sense transcripts have also been reported in mammals, insects, and fungi for genes that are part of the circadian clocks (Crosthwaite, 2004). This system coordinates the expression of some 10% of the eukaryotic genes on a daily and seasonal timescale.

### Natural ribozymes

Until about 20 years ago, it was firmly believed that proteins were the only catalytic macromolecules in biology. The discovery of the first catalytic RNA molecules, or ribozymes, in the early 1980s, however, has changed this pictures

considerably. We have already encountered several examples: RNase P, the spliceosome, and the ribosome are essentially ribozymes. In most cases, ribozymes serve an RNA-processing function using RNA as substrates. The majority of known ribozymes have been created in artificial selection experiments and hence are not a topic of this contribution; for a recent review of artificial ribozymes the interested reader is referred to Joyce (2004) and the references therein.

A number of natural ribozymes, however, are not independently stable ncRNAs but rather are part of larger RNA molecules. For example, there are four distinct groups of nucleolytic ribozymes: hammerhead and hairpin ribozymes are mostly found in plant viruses, the Varkud satellite (VS) ribozyme was found in fungal mitochondria, and hepatitis delta virus contains another ribozyme. A recent study suggests a common origin of hammerhead, hairpin, and hepatitis delta ribozymes (Harris and Elder, 2000), although convergent evolution cannot be ruled out.

The second large class of naturally occurring ribozymes is involved in the self-splicing of introns in a wide range of species; these molecules belong to one of two structural classes known as *group I* and *group II* ribozymes. All these ribozymes perform different kinds of phosphoryl transfer reactions, in which a transesterification reaction results in breakage of the backbone in the first step (Lilley, 2003). Since they behave rather like mobile genetic elements they are outside the scope of this survey; indeed many group II introns carry their own ORFs, see e.g. Zimmerly et al. (2001).

## Transcription of ncRNAs

Some ncRNAs can be found by searching for likely transcripts that do not contain an open reading frame. A survey of the *E. coli* genome for DNA regions that contain a  $\sigma$ 70 promotor within a short distance of a *Rho*-independent terminator, for instance, resulted in 144 novel possible ncRNAs (Chen et al., 2002), see also Argaman et al. (2001), Wasserman et al. (2001) and MacIntosh et al. (2001) for similar studies. This approach is limited, however, to functional RNAs that are transcribed in the “usual” manner, see Table 6. For many ncRNAs, however, the mode of transcription is unknown.

RNA Polymerase I transcribes rDNA transcription units to 18S, 5.8S and 28S rRNAs in the nucleolus. The rDNA promoters consist of the start site proximal core promoter (CP) resembling a TATA-box and an upstream control element (UCE). Both CP and UCE show poor sequence but strong structure conservation. Decreases in cell growth and protein production also reduce rRNA transcription; rRNA transcription activity oscillates during the cell cycle, showing maxima at S and G2 phase and is repressed during mitosis. In general, acetylation and phosphorylation of basal TFs regulate pol-I transcription. These modifications are performed e.g. by components of the MAPK pathway or tumor suppressors. For reviews see e.g. Paule and White (2000) and Grummt (2003).

**Table 6.** Major modes of transcription

RNA polymerase	Promoter	Location relative to start site	Transcript	Function
Pol I	Core element	-45 to +20	pre-rRNA (28S, 18S, 5.8S)	Components of the ribosome; translation
	UCE (upstream control element)	-180 to -107		
Pol II	TATA-Box Initiator	-25 to -35	mRNA	Protein coding genes
	CpG islands	-100	snRNA (U1-4)	Components of the spliceosome; mRNA splicing
			LINEs	Retrotransposon
Pol III	A-box,	+50 to +80	5S rRNA	Component of large ribosomal subunit
	B-box, C-box			
			tRNA	Translation
			snRNA (U6)	Components of the spliceosome; mRNA splicing
			7SL RNA	Component of the SRP (signal recognition particle); protein transport to ER (endoplasmic reticulum)
			SINEs	Retrotransposon

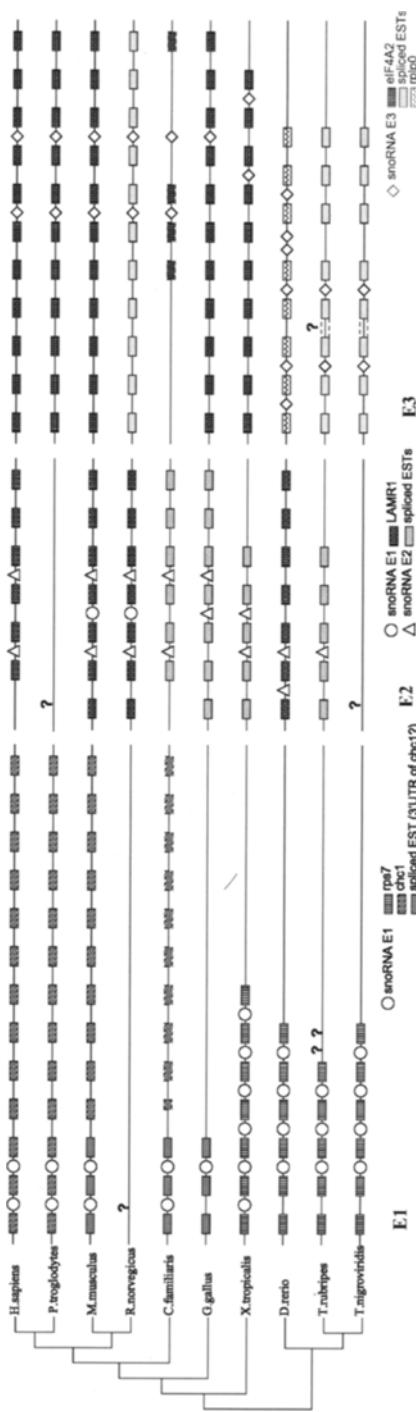
Another important class of ncRNA genes are transcribed by pol-III. Besides all canonical tRNAs and the 5S rRNA, this group includes the U6, and presumably U6acac snRNAs, RNase P and RNase MRP RNA, 7SK RNA, selenocystein tRNA, Y-RNAs, and vault RNAs (Schramm and Hernandez, 2002). Furthermore certain repetitive elements including SINEs are pol-III transcripts. For a detailed description we refer to White (1998).

The majority of transcripts is produced by pol-II, however. Most vertebrate snoRNAs are processed from introns of either protein coding genes or of “host genes” whose only known purpose is to carry an snoRNA in its intron(s) (Tycowski and Steitz, 2001), see also Bachellerie et al. (2002). Some snoRNAs, however, are transcribed directly from mono-cistronic or poly-cistronic genes, notably the U3, U8, and U13 snoRNAs. These share their promoter structure with a group of ncRNAs that contains the spliceosomal RNAs and the U7 RNA which is involved in histone mRNA processing (Hernandez, 2001). In vertebrates almost all snoRNAs are encoded in introns of a specific subclass of pol-II transcripts, the TOP genes, whose promoter elements determine a specific ratio of snoRNA and mRNA production (de Turris et al., 2004). Many vertebrate snoRNAs appear in multiple copies in different introns of the same gene, sometimes paralogs are located even on different chromosomes, see Fig. 10. The recent discovery of H/ACA snoRNA clusters within individual introns in *Drosophila* a different expression strategy for a box H/ACA snoRNA compared to box C/D snoRNAs in this species (Huang et al., 2004).

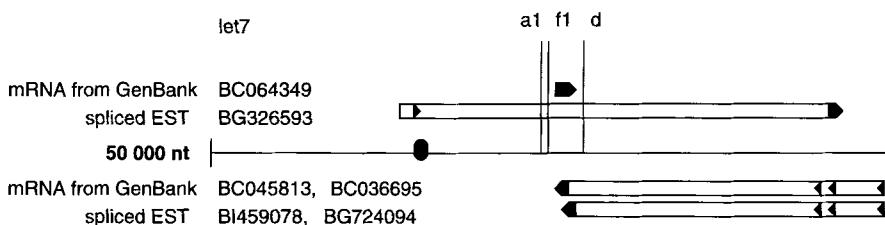
The association of intron-encoded snoRNAs with their surrounding gene surprisingly is not stable over long time-scales. U17, for example, is located in Rps7 in tetrapoda, while it is associated with the unrelated CHC1 protein in teleost fishes, Fig. 10. In addition, rodents have an additional copy of E1 in intron 3 of the lamin gene, which carries the E2 snoRNA in vertebrates. E3 switches from the ribosomal protein RPL0 to a “host gene” whose exons do not code for a functional protein.

miRNAs are processed from long primary precursors (pre-miRNAs) (Lee et al., 2002, 2003). Unlike the majority of snoRNAs, neither the genomic location of miRNAs coincides with a specific genomic context, nor is their transcription performed by a single typical mechanism. A recent survey of mammalian genomes showed that there are five major classes (Rodriguez et al., 2004): About 30% are directly transcribed by RNA Polymerase II and a 5' cap as well as a poly(A) tail is added (Cai et al., 2004), as shown for the *mir-23* cluster (Lee et al., 2004). 40% of the mammalian miRNAs are probably processed from introns (Ying and Lin, 2004, 2005) of protein coding genes, 10% of the known miRNAs reside in introns and another 13% in exons of non-coding transcripts. Antisense transcripts account for 14% of all mammalian miRNAs. The remaining cases are of uncertain transcriptional origin (Rodriguez et al., 2004).

Interestingly, our reconstruction of the duplication history of *let-7* in Fig. 7 shows that different mammalian members of this family occurs in introns of both protein-coding and non-protein-coding genes as well as in exons. Preliminary data (Tanzer et al., 2005) suggest that miRNA genes, like the genes of the three snoRNAs in



**Fig. 10.** Organization of snoRNA paralogs in the introns of their associated genes. All SnoRNAs E1, E2 and E3 identified so far, reside within introns. E1 changed its host gene from *rps7* (ribosomal protein S7) to *chc1* (chromosome condensation 1); E3 switched from *rplp0* (ribosomal protein, large, P0) to *eif4A2* (E74-like factor 4, an ets-domain transcription factor). Only snoRNA E2 remained stable associated with its host gene *LAMR1* (laminin receptor 1 [ribosomal protein SA]), which gained an additional copy of E1 in rodents. Question marks indicate incomplete genome data. For details like gene accession numbers or genome coordinates we refer to the supplemental material.



**Fig. 11.** Genomic environment of the human *let-7* family members *let-7-a1*, *let-7-f1*, and *let-7-d*. Known transcription units on the plus strand are shown above the line, the area below the line implies a location on the minus strand. The dot indicates a cluster of phylogenetic footprints (detected by tracker) that is conserved at least among amniotes. An ortholog of this particular *let-7* miRNA cluster was not found in the unfinished genome of the frog *Xenopus tropicalis* (both v2.0 and v3.0).

Fig. 10, can “move around” in the genome: Only one fourth of the known human miRNAs are located in annotated genes with known homologs in mouse and chicken. Of these 56 human miRNAs, however, less than half have known orthologs (in the Rfam miRNA Registry 4.0 which contains the results of a survey of the chicken genome) in the homologous genes.

We searched the genomic vicinity of *let-7-d*, which also contains *let-7-a1* and *let-7-f1*, for conserved non-coding DNA using the phylogenetic footprinting tool tracker (Prohaska et al., 2004). This cluster of 3 *let-7* miRNAs appears to be located within an intron of a single transcription unit (suggested by spliced EST data) with an approximate start site about 10 kb upstream of *let-7-a1*, Fig. 11. An other transcript on the plus strand with its start site between *let-7-f1* and *let-7-d* is also a possible host transcript for *let-7d*. The transcription unit OTTHUMG00000020259 (Vega data base),<sup>4</sup> which is implicated in Rodriguez et al. (2004) as carrying *let-7d* in its intron, is located on the opposite strand, however.

Approximately 500 nucleotides downstream of the transcription start site, we detected a large phylogenetic footprint cluster (100–150 nt) that is conserved among amniotes, while it is not conserved in the most closely related *let-7* cluster in Actinopterygian fishes (consisting of orthologs of *let-7-f2* and *mir-98*).

The footprint cluster does not correspond to a additional miRNA or another unannotated ncRNA in the cluster, since RNAZ (Washietl et al., 2005) classifies this region unambiguously as not containing a conserved RNA structure. A search for transcription factor binding sites within the footprint cluster using tfsearch revealed a set of common sites CREB, MZF1, GATA-1/2, Nkx-2, NrF-2 or c-Ets and Elk-1, suggesting a function in the regulation of this *let-7* miRNA family.

<sup>4</sup><http://vega.sanger.ac.uk/>

## Modifications of ncRNAs

Many, if not most, of the ncRNAs are post-transcriptionally modified. We have already encountered the snoRNA-guided pseudo-uridinylation and ribose 2'-O-methylation of rRNAs and spliceosomal RNAs. The target sites of these modification are well conserved over long evolutionary time-scales, a fact that allowed the usage of yeast rRNA methylation sites in the search for snoRNAs that modify homologous position of fruitfly rRNAs (Accardo et al., 2004).

More than 80 different nucleotide modifications are listed in the “Compilation of tRNA sequences and sequences of tRNA genes”<sup>5</sup> of various organisms (Sprinzel et al., 1996, 1998; Sprinzel and Vassilenko, 2005). These are achieved by a large family of often highly conserved enzymes, see Maas and Rich (2000) and Hopper and Phizicky (2003) for reviews. In Archaea these modifications have been shown to require four snoRNAs, one of them encoded within the intron of tRNATrp (Clouet d'Orval et al., 2001). In *S. cerevisiae* tRNA genes were shown to co-locate with the nucleolus (Thompson et al., 2003). Within this nuclear structure rRNA transcription and processing including modification by snoRNAs takes places. Recently, a selenocystein tRNA was co-immunoprecipitated with Cbf5p from Euglena, a putative pseudouridine synthase usually associated with H/ACA snoRNAs modifying rRNAs, (Russell et al., 2004). However, no tRNA modifying snoRNAs have been detected in eukaryotes so far. Studies of the evolutionary aspects of RNA editing have focused on the enzymes. Maas and Rich (2000), for instance, describe the evolution of the superfamily of RNA-dependent deaminases. We are not aware, however, of a systematic study of the evolution of the chemical modifications themselves. The 5' part of tRNAs is edited in some organisms by replacing mismatched nucleotides with nucleotides capable of forming Watson–Crick pairs in order to obtain a canonical terminal stem. This mechanism is known e.g. in the rhizopod amoeba *Acanthamoeba castellanii* and the chytridiomycete fungus *Spizellomyces punctatus* and appears to have arisen independently at least twice (Laforest et al., 2004).

In *C. elegans*, the pre-miRNA of *let-7* has been shown to undergo *trans*-splicing to the spliced leader 1 (SL1) RNA. This process allows folding of the pre-miRNA such that the miRNA precursor (pre-miRNA) forms a stem-loop structure, which in turn is cleaved by nuclear RNase III Drosha (Filippov et al., 2000; Carmell and Hannon, 2004; Lee et al., 2002, 2003). Both, the mature miRNAs and the pre-miRNA can undergo A to I RNA editing by an RNA-specific adenosine deaminase (ADAR) (Luciano et al., 2004).

## RNA motifs associated with protein-coding mRNAs

### mRNA structure

In contrast to ncRNAs, the primary function of mRNAs is to encode in its exons the information that allows the translation machinery to generate proteins. Exon

<sup>5</sup><http://www.tRNA.uni-bayreuth.de/>

recognition by the spliceosome can be affected by many features of the pre-mRNA including exon length, promoter architecture, the presence of enhancer and silencer elements, the strength of splicing signals, and RNA processivity. It has also been proposed repeatedly that pre-mRNA secondary structures influence splicing activity. A recent review of the topic (Buratti and Baralle, 2004) strongly suggests that many pre-mRNA sequences contain selected regions folding *in vivo* into well-defined secondary structures that are likely to play a role in the splicing process.

Eukaryotic mature mRNA exhibit a tripartite structure: an untranslated region at the 5' end (5'UTR), the coding regions which is translated into amino acids, and an untranslated region at the 3' end (3'UTR). With the exception of the replication dependent histone genes mentioned above, the 3' end of the mature mRNA carries a poly(A) tail. Both untranslated regions are involved in the post-transcriptional regulation of gene expression processes, like subcellular localization, mRNA stability and translation efficiency (Le and Maizel, 1997; Day and Tuite, 1998; Oleynikov and Singer, 1998; Pesole et al., 2001; Wilkie et al., 2003). These processes are mainly controlled by *cis*-acting functional elements in the UTRs, which comprise both sequence motifs and RNA structure motifs. Short sequence motifs may be potential binding sites for *trans*-acting factors, while longer sequence motifs found in UTRs have been hypothesized to be antisense RNA binding sites (Lipman, 1997). In addition, a number of motifs are known that are determined by structural features rather than nucleic acid sequence.

In general, the protein coding region of mRNAs is much better conserved than the UTRs. The distribution of conserved sequence motifs in the UTRs is not uniform: the 3'UTR is typically better conserved than 5'UTR and introns (Jareborg et al., 1999). For example, 30% of 3'UTRs in different vertebrate mRNAs contain highly conserved regions which are at least 100 nt long and show at least 70% similarity (Duret et al., 1993). The overall higher conservation of the 3'UTR may be a consequence of the observation that post-transcriptional regulation in 3'UTRs is rather based on protein complexes than on single proteins (Wilkie et al., 2003). Antisense binding, as e.g., in the case of miRNAs, also leads us to expect many non-structural binding motifs. This is in particular the case in plants, where miRNA targets typically closely match the corresponding miRNA, see Kidner and Martienssen (2005) and the references therein. (In mammals, however, the requirements for mRNA-miRNA interactions appear to be much more complex (Vella et al., 2004).) In contrast, 5'UTRs regulatory motifs might be mostly structural motifs. It is known, for example, that translation initiation is essentially controlled by RNA structures in 5'UTRs (Day and Tuite, 1998; Pesole et al., 2001).

A more detailed analysis showed, however, that pattern of conservedness is reversed at the border of the coding region. The 30 nt of the 5'UTR immediately upstream of the start codon is the best conserved regions with 70–80% sequence identity between human and mouse; in contrast, the 3'UTR is very poorly conserved immediately downstream of the stop codon (Miyata et al., 1980; Shabalina et al., 2004). This pattern can be explained by the specific interaction with sequence specific binding factors initiating translation in the 5'UTR on the one hand, and the fact that

the first segment of the 3'UTR is covered by ribosome and hence inaccessible to specific factors at the termination of translation.

The UTRs and the coding regions are subject to different functional constraints and hence evolve differently; even the 5'UTRs and 3'UTRs of the same gene do not necessarily share the same evolutionary dynamics (Lariza et al., 2002). In addition there are also mRNAs that encode nearly identical proteins but have highly diverged UTRs (Duret et al., 1993; Jareborg et al., 1999; Lariza et al., 2002), suggesting that the divergent UTRs form specific translational regulation patterns which enables them to reply differently to variable stimuli.

### Detection of UTR motifs

A handful of *cis*-acting regulatory motifs in mRNAs have been characterized experimentally; these are collected in the UTRsite (Pesole et al., 2002) and Transterm data bases (Dalphin et al., 1999). Functional RNA structures in UTRs are in general not as long as ncRNAs, since they are limited by the size of their UTRs. The average length of human UTRs is about 210 for 5'UTRs and about 1027 for 3'UTRs (Pesole et al., 2001). RNA structure motifs in UTRs can thus be expected to be relatively small, simple structures. This limits the usable information, and hence the fast and reliable prediction of structural regulatory elements in UTRs has remained a largely unsolved problem.

Standard sequence alignment procedures usually fail to align UTRs in a meaningful way (Jareborg et al., 1999). Detecting structural motifs in UTRs will therefore require algorithms that optimize sequence alignment and secondary structure simultaneously. Existing methods, which characterize putative motifs automatically, can be classified in (a) methods which require the description of a motif and search for similar instances of such a motif or (b) methods which search for motifs that are significantly overrepresented in a dataset.

While most regulatory motifs found in UTRs are conserved in secondary structure, some known motifs show conservation on the sequence level. In general such sequence motifs do not require an exact nucleotide substring, but allow some variation in nucleotide composition, or may consist of several conserved fragments separated by unconserved regions. Fragmented motifs may, for example, occur in regulatory structures, which exhibit sequence conservation in loops.

Many of the tools developed to identify functional RNA motifs in general have also been applied to UTRs. Among these Palingol (Billoud et al., 1996), PatSearch (Pesole et al., 2000b) and RNAMotif (Macke et al., 2001) identify instances of a previously defined motif descriptor. Because of the limited information available in a UTR motif it is hard to define descriptors that are both specific and sensitive by hand. Other tools, therefore, are designed to require only limited information about a known motif and recognize motif features, which discriminate sequences containing the motif from sequences not containing it, automatically (see ERPIN, Gautheret and Lambert (2001)).

An even harder problem has to be solved when the motif is completely unknown. The detection problem in such a case can be treated as a classification problem:

From an arbitrarily given set of UTRs all UTRs sharing a common motif shall be classified in the same group. The first approaches to this problem are implemented in the tools comRNA (Ji et al., 2004), which identifies novel mRNA structure motifs by clustering similar stems, and RNAPProfile (Pavesi et al., 2004), which identifies the most conserved motif in a set of sequences where at least some share the same common motif.

### Important regulatory motifs in UTRs of mRNAs

Gene expression is controlled by *cis*- and *trans*-acting factors during both transcription and during translation. By regulating translation a cell is able to respond quickly to environmental changes. The mature mRNA already resides in the cellular plasm but the amount and type of protein which will be translated depends on several cellular conditions. *Cis*-acting elements in the untranslated regions of mature mRNA bind *trans*-acting factors and control in this way translational efficiency, mRNA stability and subcellular localization. A selection of examples of such regulatory motifs in UTRs will be given here, see also Table 7.

**Iron response elements (IRE):** Are short hairpin structures with an internal loop and a conserved sequence in the hairpin loop, which are observed in 5' UTRs of ferritin mRNAs in 3'UTRs of transferrin receptor mRNAs (Hentze and Kühn, 1996). They can be classified in two slightly different instances, the first containing an internal loop of length three, which is replaced by a bulge loop in the second. Both have the primary consensus motif CNNNNNCAGWGH (Pesole et al., 2002). The IRE motif can be readily described with regular grammars; because of the highly redundant sequence pattern and frequent, simple secondary structure one has to expect a large number of false positives, however.

**Translation control elements (TCE)** are short elements (~90 nt) found in the 3'UTR of nanos mRNA of drosophila (Crucks et al., 2000). Its secondary structure is composed of a helix and a multiloop with two hairpin loops branching off, one with a conserved sequence in the hairpin loop.

**Internal ribosome entry site (IRES) elements:** were first described in the 5'-untranslated region of picornavirus RNA (Jang et al., 1988). The IRES element enables cap-independent initiation of translation starting at an internal initiation codon. In addition to several types of viruses, which contain an IRES element, a small group of eukaryotic mRNA can be translated by internal ribosome entry. IRES-containing mRNAs mostly encode regulatory proteins such as, e.g., growth factors and transcription factors. Several studies have reported that under stress conditions, where cap-dependent translation is blocked, translation of specific mRNAs is enabled through IRES elements (Martineau et al., 2004 and references therein). Another function of IRESs involves the control of alternative initiation of translation. For example, the human fibroblast growth factor 2 contains 5 translation initiation codons. Translation initiation of the codon proximal to the 5'-end is initiated by a cap-dependent process, whereas initiation of the remaining codons depends on the IRES (Bonnal et al., 2003). IRES elements are defined by functional criteria and cannot yet be predicted by the presence of characteristic RNA

sequence or structural motifs. In general, there are no significant similarities between individual IRESs unless they are from related sources.

*Selenocysteine insertion sequences* (SECIS): Is found in the coding region of some eubacterial mRNAs and in 3' untranslated regions of some mRNAs in archaea and eukaryotes (Krol, 2002). In eubacteria, it forms a hairpin structure of conserved length with the selenocysteine codon in the outer helix. In archaea, the primary rather than the secondary structure is conserved. The consensus is a hairpin structure that differs in stem length, occurrence of internal loops and size of the hairpin loop, but it has a very conserved sequence motif in the helix beneath the apical loop. In eukaryotes, the secondary structure contains most of the information while only small sequence motifs are conserved. The core secondary structure is composed of a long hairpin structure consisting of two (type 1) or three (type 2) consecutive helices (Fagegaltier et al., 2000; Krol, 2002).

At present, it is unclear whether large regulatory motifs such as IRES, IRE, or SECIS elements, arose independently in different genes or gene families or whether there are mechanisms that allow their lateral spread within a genome.

*Novel UTR Motifs:* In addition to post-transcriptional regulatory mechanisms that are specific to a particular gene or gene family, there exist also mechanisms which are observed in a broader range of mRNAs (Table 7). Such relatively non-specific regulatory processes are characterized by similar primary and/or secondary structures in mRNAs of different genes of the same organism. We performed a search for sequence elements of this type in the human genome.

Using NCBI blast (Altschul et al., 1997), we compute pairwise alignments of repeat masked human UTRs from Ensembl data base (release 24). The majority of UTRs were not conserved on the sequence level, suggesting that also non-specific regulatory motifs show large sequence divergence. Furthermore, we found many more conserved sequence blocks in 3'UTRs than in 5'UTRs. From the pairwise alignments we built a weighted similarity graph to identify clusters of UTRs with conserved regions by complete linkage clustering (Day and Edelsbrunner, 1984). Alternative transcripts were not allowed to occur in the same cluster. Sequences of each cluster were aligned using dialign2 (Morgenstern, 1999) in order to identify putative regulatory sequence motifs.

We then used RNAz (Washietl et al., 2005) to check whether some of these multiple alignments contain conserved RNA secondary structures. Among 481 5'UTR multiple alignments, 10% had regions forming with high probability stable RNA structures. Among the set of 1223 multiple 3'UTRs alignments, 21% alignments contain stable RNA structures. Table 8 lists the annotation of the best RNA predictions using infernal and the Rfam data base. All significant hits matched the iron response element. The corresponding genes, however, are not known to be involved in the iron metabolism. We suspect that at least some of these cases form IRE-like structures that do not function as IREs, indicating that still more specific descriptors for UTR elements including IRE are desirable.

The small fraction of RNA motifs with known function that was recovered in our survey suggests that most non-gene-specific mRNA motifs have very little well-conserved sequence information and most of them, including IRES, SECIS, IRE,

**Table 7.** Important general regulatory motifs found in mature eukaryotic mRNAs (Pesole et al., 2001)

Motif	Function	Description	References
5'UTR <i>m<sup>7</sup>G</i> cap structure	Stabilization, initiation	Prevents processing of mRNA from 5' to 3' end, hence stabilizes mRNA; eIFs bind to cap, which governs pre-initiation complex with small ribosomal subunit and initiates scanning	Day and Tuite (1998) and Gebauer and Hentze (2004)
Initiation codon	Initiation, translational efficiency	Efficiency of translation start recognition depends on primary sequence context of AUGcodon; optimal context for vertebrates is (A/G) CCAU GG	Kozak (1989), Day and Tuite (1998) and Pesole et al. (2000a)
uORF	Translational efficiency	Inhibits translation by leaky scanning: scanning complex may either bypass upstream start codon depending on sequence context and mean ORF is translated or may start translation at upstream start codon	Day and Tuite (1998) and Gebauer and Hentze (2004)
IRES	Initiation	Alternative to ribosomal scanning; pre-initiation complex interacts with IRES element and scanning starts at this site	Le and Maizel (1997), Kozak (2001) and Mignone et al. (2002)
Stable RNA structures	Translational efficiency	Very stable secondary structures in 5'UTRs can impede scanning	Kozak (1991), Mignone et al. (2002), Wilkie et al. (2003) and Gebauer and Hentze (2004)
Repeats	Initiation, translational efficiency	<i>Ahu</i> -elements in 5'UTRs e.g. repress translational efficiency; reason may be repression of initiation by <i>Ahu</i> -elements forming stable secondary structures or containing weak start codons	Landy et al. (2001)

3'UTR Zipcodes	Localization	RNA binding proteins bind to different zip codes and direct mRNA to subcellular region where corresponding protein is translated; proteins recognize zipcode by primary and tertiary structure	Oleynikov and Singer (1998) and Hesketh (2004)
Poly(A)-tail	Stabilization, initiation	Prevents degradation of mRNA target sites	Day and Tuite (1998), Wilkie et al. (2003) and Gebauer and Hentze (2004)
CPE	Stabilization, translational efficiency	CPEB binds at CPE and induces polyadenylation; a complex of CPEB and Maskin bound to CPE interacts with cap structure by binding to the eIF4F complex and translation is repressed	Mendez and Richter (2001), Wilkie et al. (2003) and Gebauer and Hentze (2004)
AREs	Stabilization	Influence rate of deadenylation depending on type of AU-rich element (ARE 1, ARE 2 or ARE 3) mRNAs evoke degradation of mRNA by imperfect base-pairing interactions with 3'UTR	Mignone et al. (2002) and Meisner et al. (2004)
mRNA target sites	Destabilization	CAG/CUG repeats in 3'UTRs e.g. result in very long mRNAs, which show in yeast different subcellular distribution	Wilkie et al. (2003) and Gebauer and Hentze (2004)
Repeats	For e.g. localization		Fabre et al. (2002) and Mignone et al. (2002)

Most regulatory elements influence initiation of translation. Regulatory elements specific to mRNAs of particular genes are not listed. **Abbreviations:** IRES = internal ribosome entry site, AREs = AU-rich elements, uORF = upstream open reading frame, CPE = cytoplasmic polyadenylation element, CPEB = cytoplasmic polyadenylation element binding protein, ACE = adenylylate control element, eIFs = eukaryotic initiation factors, PABP = poly(A)-binding protein, UTR = untranslated region. Repeats found in UTRs include short interspersed elements (SINEs), long interspersed elements (LINEs), mini- and micro-satellites (Mignone et al. (2002)); these are not listed above.

**Table 8.** RNA structure annotation based on infernal and description of corresponding genes

Rfam model	Score	Genes	Description (ensembl release 24)
<b>5' UTR</b>			
IRE	6.8	ENSG00000120853	—
		ENSG00000166104	—
IRE	6.31	ENSG00000092199	Heterogeneous nuclear ribonucleoprotein
		ENSG00000159267	Biotin-protein ligase
IRE	9.89	ENSG00000129873	Testis-specific chromodomain protein Y2
		ENSG00000172288	Testis-specific chromodomain protein Y1
		ENSG00000172353	Testis-specific chromodomain protein Y1
<b>3' UTR</b>			
IRE	8.24	ENSG0000066294	CD84 antigen
		ENSG00000134822	Fatty acid desaturase
IRE	8.26	ENSG00000165282	Phosphatidylinositol-glycan biosynthesis
		ENSG00000152056	Sigma-adaptin 1C
IRE	8.28	ENSG00000110436	Amino acid transporter 2
REN-SE	12.1	ENSG00000171596	G protein-coupled receptor 66
		ENSG00000166676	—
		ENSG00000090659	CD209 antigen
IRE	8.33	ENSG00000181894	—
		ENSG00000149451	ADAM 33 precursor
IRE	11.29	ENSG00000185753	—
		ENSG00000064115	Transmembrane 7 superfamily protein member 3 precursor
IRE	12.51	ENSG0000012048	Breast cancer type 1 susceptibility protein
		ENSG00000156675	Rab coupling protein
		ENSG00000142687	Polycystic kidney disease 1-related
IRE SECIS	9.19	ENSG00000181719	—
	10.87	ENSG00000178887	—
		ENSG00000180747	—

and many others depend crucially on secondary structure. On the other hand, we detected hundreds of statistically significant sequence patterns that occur in multiple RNAs for which so far no function has been described. A pattern is defined by the consensus sequence of a run of gapless columns in the multiple alignments.

One possible function of sequence patterns in the 3'UTR of mRNAs is to act as target sites for miRNAs (Wilkie et al., 2003; Gebauer and Hentze, 2004). We therefore tested all gapless regions in the multiple alignments for potential miRNA target sites. To this end we used the collection of all human miRNAs from Rfam

(release 5.0, September 2004) (Griffiths-Jones, 2004) and two different miRNA target prediction tools: miRanda (Enright et al., 2003) and RNAhybrid (Rehmsmeier et al., 2004). Table 9 lists the best-scoring candidates. In contrast, an analysis of the 5'UTRs and their flanking regions did not yield a potential site located within the untranslated region of mRNA that was predicted by both methods.

### RNA structures in coding regions

It is widely believed that RNA structures in ORFs can interfere with translation, although this phenomenon has not been studied systematically to our knowledge (Katz and Burge, 2003). It is plausible to assume that coding regions are therefore largely devoid of secondary structures. There are, however, a number of well-known exceptions to this rule. A variety of conserved secondary structure elements have been detected in computational surveys of single stranded RNA virus genomes (Hofacker et al., 2004b; Tuplin et al., 2002, 2004; Thurner et al., 2004; Witwer et al., 2001). A comparative study of 28 different species (Katz and Burge, 2003) provides evidence for wide-spread selection for local secondary structures in mRNAs, in particular in eubacteria. Most recently, (Pedersen et al., 2004a, b) devised an SCFG-based algorithm for detecting conserved secondary structures motifs specifically within coding sequences.

The *Rev Response Element* (RRE), for example, forms a five-fingered motif spanning some 300 nt (Dayton et al., 1992), located in the *env* gene of HIV. The structure is well conserved among diverse HIV strains, see, e.g., Hofacker et al. (1998), Hofacker and Stadler (1999) and Konecny et al. (2000). The interaction of RRE with the *Rev* protein reduces splicing and increases the transport of unspliced and single-spliced transcripts to the cytoplasm, which is necessary for the formation of new virion particles (Malim et al., 1989).

A *cis-acting regulation element* (CRE) within the coding region of several picornaviruses has been described in a number of different picornaviruses. The function of the CRE probably involves the initiation of the synthesis of the negative-sense strand template RNA during virus replication (Goodfellow et al., 2003). The CRE has been found as in a computational survey (Witwer et al., 2001) in most genera of the picornaviridae. Interestingly, its genomic location varies between genera.

The best known example in a higher organism is the stem-loop structure in the coding region of the *ASH1* gene of yeast which localizes the *ASH1* mRNA to the bud tip (Chartrand et al., 1999). With the exception of a viral elements, however, the functions, as well as possible evolutionary relationships, of structured RNA motifs within ORFs remain unknown.

### Riboswitches

Some RNA molecules exhibit two competing conformations, whose equilibrium can be shifted easily by molecular events such as the binding of another molecule.

**Table 9.** Potential miRNA target sites in human 3' untranslated regions

miRNA	Gene	Score				Protein family	
		ENSG00000...	Protein	miRand	RNAhybrid	ENSF0000000...	Name
hsa-miR-187	183850	Zinc finger protein 254	—	11.65	0.000096	0001	Zinc Finger
hsa-miR-187	181342	—	AKAP-binding sperm protein ropporin	11.11	0.000081	0001	Zinc Finger
hsa-miR-134	065371	—	AKAP-binding sperm protein ropporin	10.50	0.000914	5520	—
hsa-miR-134	114547	—	AKAP-binding sperm protein ropporin	10.50	0.000914	5520	—
hsa-miR-324-5p	129277	—	Small inducible cytokine A4 precursor	8.41	0.000062	0592	—
hsa-miR-324-5p	189315	—	Small inducible cytokine A4 precursor like	8.36	0.000062	0592	—
hsa-miR-184	177111	—	—	11.49	0.000143	2097	DPY19
hsa-miR-184	177990	—	—	11.03	0.000143	2097	DPY19

We report z-scores for miRand and p-values for RNAhybrid as computed by these tools. Gene annotation is taken from Ensembl (release 24).

This can be used to regulate gene expression, when the two mutually exclusive alternatives correspond to an active and in-active conformation of the transcript (Merino and Yanofsky, 2002). Mechanistically, one fold of the mRNA, the repressing conformation, contains a terminator hairpin or some other structural element which conceals the translation initiation site, whereas in the alternative conformation, the non-repressing one, the gene can be expressed (Henkin and Yanofsky, 2002). An early computational study concluded that RNA switches are readily accessible in evolution and are therefore probably not exceptional instances of unusual RNA behavior (Flamm et al., 2000). The use of two competing RNA conformations allows molecular events like the binding of a target metabolite by a protein to influence which of the alternative conformations the terminator or the anti-terminator is formed, hence coupling the gene expression to the concentration of the target metabolite.

The best known example of such behavior are the riboswitches (Vitreschak et al., 2004). These are autonomous structural elements primarily found within the 5'-UTRs of bacterial mRNAs, which, upon direct binding of small organic molecules, can trigger conformational changes, leading to an alteration of the expression for the downstream located gene. Their general architecture shows two modular units (Winkler and Breaker, 2003), a ligand-binding one, which function as a “sensor” for a small metabolite and a unit which “interprets” the signal from the “sensor” unit and interfaces to those RNA elements involved in gene expression regulation. The size of the “sensor”-unit ranges typically from 70–170 nt, which is unexpectedly large compared to artificial aptamers obtained by *in vitro* directed evolutionary experiments. While for most riboswitches the ligand-binding domain is highly conserved among various organisms, the “interpretation” module varies strongly in sequence, structure and mechanism by which it controls the appended gene. Riboswitches and engineered allosteric ribozymes (Breaker, 2002; Silverman, 2003) demonstrate impressively that RNA is indeed capable of maintaining a complex metabolic state without the help of proteins.

Riboswitches regulate several key metabolic pathways (Brantl, 2004; Nudler and Mironov, 2004) in bacteria including those for coenzyme B<sub>12</sub>, thiamine, pyrophosphate, flavin monophosphate, S-adenosylmethionine and a couple of important amino acids. The search for additional elements is ongoing, e.g. Barrick et al. (2004) and Lesnik et al. (2005). The program Riboswitch finder (Bengert and Dandekar, 2004) utilizes consensus motifs of known elements to detect new prokaryotes riboswitches.

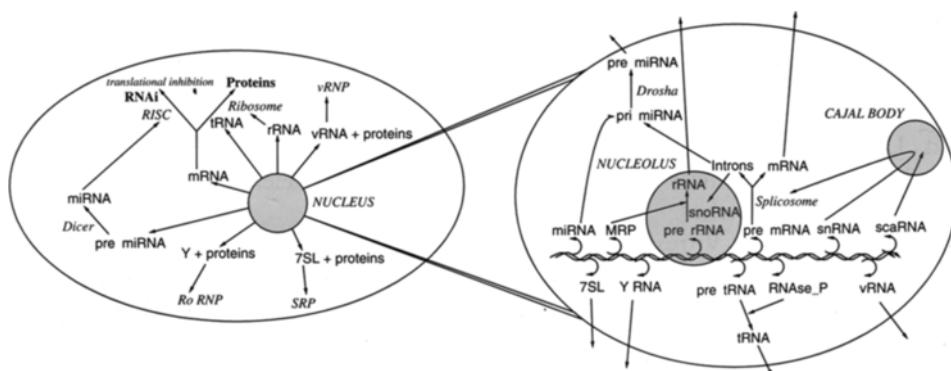
A recent paper by Vitreschak et al. (2003) applied comparative and phylogenetic analysis to vitamin B<sub>12</sub>-related genes using 200 sequences from 66 bacterial genomes. They identified a highly conserved regulatory RNA structure, the B12-element a cobalamin riboswitch, which is widely distributed in 5'-UTRs of vitamin B<sub>12</sub>-related genes in eubacteria. Comparison of the reconstructed phylogenetic tree for the B12-element with standard trees showed both lineage- and gene-specific branches, as well as a large number of recent gene duplications and horizontal gene transfer events. A related study is reported in (Nahvi et al., 2004). Comparative approaches were also used to study the L-box regulon regulating the lysine synthesis pathway (Grundy

et al., 2003) and the S- and T-boxes in the methionine metabolism of Gram-positive bacteria (Rodionov et al., 2004). While most riboswitches were found in bacteria, such metabolite-binding RNA domains are also present in some eukaryotic genes (Sudarsan et al., 2003). These findings, and the fact that riboswitches bind their effectors directly without the need of additional factors, suggest that riboswitches represent one of the oldest regulatory systems.

## Concluding remarks

The recent discoveries in the “modern RNA World” have made it obvious that large-scale mRNA expression profiling data can provide only a partial picture of gene expression. Most post-transcriptional events are mediated by the association of RNAs with specific proteins or macromolecular protein complexes. Comprehensive determination of the RNA targets of RNA-binding proteins is therefore likely to be important in deciphering the complex events at this level of gene regulation. Approaches to exploring the post-transcriptional RNA world with DNA micro-arrays are discussed e.g. in Iyer (2004).

Fig. 12 gives a sketch of the probably most ancient part of the RNA-based regulation system of the eukaryotic cell: ncRNAs in Eukaryotic cells seem to fall into two major groups according to their subcellular localization and thus function. The nuclear fraction mainly performs ncRNA processing and maturation. SnRNAs, snoRNAs, and scaRNAs seem to be the major players, forming the central part of the nuclear RNA regulatory network. They modify themselves as well as other ncRNAs including rRNAs, but maybe even tRNAs. Another RNA processing mechanism, RNAediting, in general does not require guide RNAs as in the case of kinetoplasts. Besides, the snRNAs act on coding hnRNA (pre-mRNA) by splicing



**Fig. 12.** Subcellular localization of ncRNAs in eukaryotic cells. The figure contains all eukaryotic ncRNA mentioned in this publication. The nuclear fraction of these regulators (right) functions in processing of ncRNAs, whereas the cytoplasmic ones (left) are involved in translation.

introns. Upon export to the cytoplasm, the majority of ncRNAs is involved in protein translation. miRNAs regulate protein expression by translational inhibition or RNAi, 7SL RNA transports mRNA of secretory proteins to the ER (endoplasmatic reticulum). Coding and ncRNAs thus share a similar “life cycle” depending on their subcellular localization: regulation and maturation is performed in the nucleus, their work is done in the cytoplasm.

It is commonly assumed that the primordial cell looked much more like a bacterial than a eukaryotic cell. For a discussion of the origin of the eukaryotic cell and its mitochondria we refer the reader to Lang et al. (1999) and Andersson et al. (2003). Because of the lack of a nuclear membrane, transport mechanisms were not required. Furthermore, intronless genomes did not require a splicing-like mechanism. Instead, poly-cistronic transcripts of ncRNA and/or mRNA might have been processed by RNA modification and subsequent endonucleolytic cleavage. This picture is consistent with our present knowledge of the evolutionary history of the major ncRNA families summarized in Fig. 13.

While some RNAs, in particular those involved in protein synthesis predate the Last Universal Common Ancestor of all extant life forms, novel RNA families with novel—mostly regulatory—function have been invented throughout the history of life. The picture in Fig. 13 is almost certainly incomplete due to a bias in the available data which are concentrated on a small number of well-studied model organisms (mainly vertebrates, arthropods, nematods, yeast, rice, arabidopsis and bacteria). The recent discovery of a novel class of expressed ncRNAs with unknown function in *D. discoideum* (Aspegren et al., 2004) and the large number of still poorly understood bacterial sRNAs, see e.g. Hershberg et al. (2003) and Altuvia (2004),

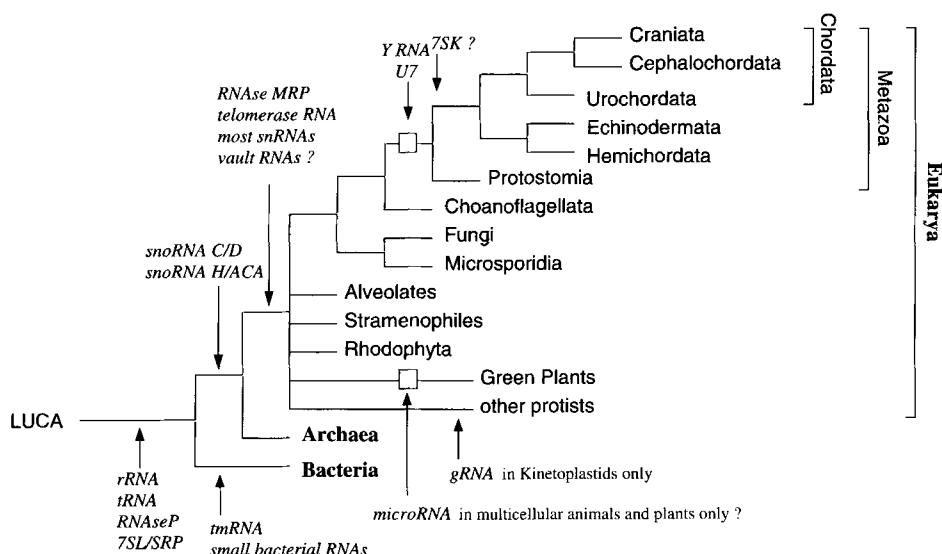


Fig. 13. Evolutionary origin of the most prominent ncRNA families.

suggests that quite a few ncRNA innovations in less-studied lineages could have escaped our attention so far.

The evidence compiled in this contribution indicates an explosive expansion of some ncRNA families, in particular of miRNAs, in the vertebrate lineage. Higher plants might show a similar pattern. In both cases, genome duplications are a plausible mechanism that at least contributed to expansion. Multiple dispersed copies of some snRNAs, in contrast, can be explained by the recent observation that certain retroviruses package and reverse-transcribe snRNAs (Giles et al., 2004). Usually, this mechanism produces pseudogenes that are associated with LTRs of endogenous retroviruses. The mechanism or mechanisms that lead to duplicates of intron-encoded snoRNAs, or the processes leading to a change from intronic to exonic expression in paralogous miRNAs, on the other hand, still remain in the dark.

We close our discussion by emphasizing that it is by no means complete: topics such as the relationships of ncRNAs and repetitive elements (e.g. Alus) or mobile genetic elements (e.g. group II introns or endogenous retroviruses) have been neglected here.

## Acknowledgements

We thank Rolf Backofen and Daniel Gautheret for their comments on an earlier version of this manuscript. This work was supported in part by the Austrian *Fonds zur Förderung der Wissenschaftlichen Forschung*, Project No. P15893, by the German DFG Bioinformatics Initiative BIZ-6/1-2, and by the Austrian *Gen-AU bioinformatics integration network* sponsored by BM-BWK and BM-WA.

## References

- Accardo, M.C., Giordano, E., Riccardo, S., Digilio, F.A., Iazzetti, G., Calogero, R.A., Furia, M., 2004. A computational search for box C/D snoRNA genes in the *D. melanogaster* genome. *Bioinformatics* 20, 3293–3301.
- Adai, A., Johnson, C., Mlotshwa, S., Archer-Evans, S., Manocha, V., Vance, V., Sundaresan, V., 2005. Computational prediction of miRNAs in *Arabidopsis thaliana*. *Genome Res.* 15, 78–91.
- Allen, E., Xie, Z., Gustafson, A.M., Sung, G.-H., Spatafora, J.W.S., Carrington, J.C., 2004. Evolution of microRNA genes by inverted duplication of target gene sequences in *Arabidopsis thaliana*. *Nat. Genet.* 36, 1282–1290.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J., 1990. Basic local alignment search tool. *J. Mol. Biol.* 215, 403–410.
- Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman, D.J., 1997. Gapped BLAST and PSI-BLAST: a new generation of protein data base search programs. *Nucleic Acids Res.* 25, 3389–3402.
- Altuvia, S., 2004. Regulatory small RNAs: the key to coordinating global regulatory circuits. *J. Bacteriol.* 186, 6679–6680.
- Ambros, V., Lee, R.C., Lavanway, A., Williams, P.T., Jewell, D., 2003. MicroRNAs and other tiny endogenous RNAs in *C. elegans*. *Curr. Biol.* 13, 807–818.

- Amores, A., Force, A., Yan, Y.L., Joly, L., Amemiya, C., Fritz, A., Ho, R.K., Langeland, J., Prince, V., Wang, Y.L., Westerfield, M., Ekker, M., Postlethwait, J.H., 1998. Zebrafish *Hox* clusters and vertebrate genome evolution. *Science* 282, 1711–1714.
- Andersen, A.A., Panning, B., 2003. Epigenetic gene regulation by noncoding RNAs. *Curr. Opin. Cell Biol.* 15, 281–289.
- Andersson, S.G.E., Karlberg, O., Canbäck, B., Kurland, C.G., 2003. On the origin of mitochondria: a genomics perspective. *Philos. Trans. R. Soc. London B: Biol. Sci.* 358, 165–177.
- Argaman, L., Vogel, J., Bejerano, G., Wagner, E., Margalit, H., Altuvia, S., 2001. Novel small RNA-encoding genes in the intergenic regions of *Escherichia coli*. *Curr. Biol.* 11, 941–950.
- Aspegren, A., Hinss, A., Larsson, P., Larsson, A., Söderbom, F., 2004. Novel non-coding RNAs in *Dicystostelium discoideum* and their expression during development. *Nucleic Acids Res.* 32, 4646–4656.
- Avner, P., Heard, E., 2001. X-chromosome inactivation: counting, choice, and initiation. *Nat. Rev. Genet.* 2, 59–67.
- Azzouz, T.N., Schümperli, D., 2003. Evolutionary conservation of the U7 small nuclear ribonucleoprotein in *Drosophila melanogaster*. *RNA* 9, 1532–1541.
- Babitzke, P., Yanofsky, C., 1993. Reconstitution of *Bacillus subtilis* Trp attenuation in vitro with TRAP, the Trp RNA-binding attenuation protein. *Proc. Natl. Acad. Sci. USA* 90, 133–137.
- Bachellerie, J.-P., Cavallé, J., Hüttenhofer, A., 2002. The expanding snoRNA world. *Biochimie* 84, 775–790.
- Bafna, V., Zhang, S., 2004. FastR: fast database search tool for non-coding RNA. Proceedings of the IEEE Computer and Systems Bioinformatics Conference.
- Bailey, S., Wichterlechkarn, J., Johnson, D., Reilly, B.E., Anderson, D.L., Bodley, J.W., 1990. Phylogenetic analysis and secondary structure of the *Bacillus subtilis* bacteriophage RNA required for DNA packaging. *J. Biol. Chem.* 265, 22365–22370.
- Barad, O., Meiri, E., Avniel, A., Aharonov, R., Barzilai, A., Bentwich, I., Einav, U., Gilad, S., Hurban, P., Karov, Y., Lobenhofer, E., Sharon, E., Shibolet, Y., Shtutman, M., Bentwich, Z., Einat, P., 2004. MicroRNA expression detected by oligonucleotide microarrays: system establishment and expression profiling in human tissues. *Genome Res.* 14, 2486–2494.
- Barrick, J.E., Corbino, K.A., Winkler, W.C., Nahvi, A., Mandal, M., Collins, J., Lee, M., Roth, A., Sudarsan, N., Jona, I., Wickiser, J.K., Breaker, R.R., 2004. New RNA motifs suggest an expanded scope for riboswitches in bacterial genetic control. *Proc. Natl. Acad. Sci. USA* 101, 6421–6426.
- Bartel, D.P., Chen, C.-Z., 2004. Micromanagers of gene expression: the potentially wide-spread influence of metazoan microRNAs. *Nat. Genet.* 5, 396–400.
- Baskerville, S., Bartel, D.P., 2002. A ribozyme that ligates RNA to protein. *Proc. Natl. Acad. Sci. USA* 99, 9154–9159.
- Beja, O., Ullu, E., Michaeli, S., 1993. Identification of a tRNA-like molecule that copurifies with the 7SL RNA of *Trypanosoma brucei*. *Mol. Biochem. Parasitol.* 57, 223–229.
- Ben-Shlomo, H., Levitan, A., Shay, N.E., Goncharov, I., Michaelim, S., 1999. RNA editing associated with the generation of two distinct conformations of the trypanosomatid leptomomas collosoma 7SL RNA. *J. Biol. Chem.* 274, 25642–25650.
- Bengert, P., Dandekar, T., 2004. Riboswitch finder: a tool for identification of riboswitch RNAs. *Nucleic Acids Res.* 32 (Web Server Issue), W154–W159.
- Bennasser, Y., Le, S.Y., Yeung, M.L., Jeang, K.T., 2004. HIV-1 encoded candidate micro-RNAs and their cellular targets. *Retrovirology* 1, 43 (Epub).
- Berezikov, E., Guryev, V., van de Belt, J., Wienholds, E., Plasterk, R.H.A., Cuppen, E., 2005. Phylogenetic shadowing and computational identification of human microRNA genes. *Cell* 120, 21–24.
- Berteaux, N., Lottin, S., Adriaenssens, E., Van Coppennolle, F., Leroy, X., Coll, J., Dugimont, T., Curgy, J.-J., 2004. Hormonal regulation of H19 gene expression in prostate epithelial cells. *J. Endocrinol.* 183, 69–78.
- Bililou, B., Kontic, M., Viari, A., 1996. Palingol: a declarative programming language to describe nucleic acids' secondary structures and to scan sequence data bases. *Nucleic Acids Res.* 24, 1395–1403.
- Bishop, K.N., Holmes, R.K., Sheehy, A.M., Malim, M.H., 2004. APOBEC-mediated editing of viral RNA. *Science* 305 (5684), 645.
- Blattner, F.R., Plunkett III, G., Bloch, C.A., Perna, N.T., Burland, V., Riley, M., Collado-Vides, J., et al., 1997. The complete genome sequence of *Escherichia coli* k-12. *Science* 277, 1453–1474.
- Blencowe, B.J., 2002. Transcription: surprising role for an elusive small nuclear RNA. *Curr. Biol.* 12, R147–R149.

- Bonnal, S., Schaeffer, C., Creancier, L., Clamens, S., Moine, H., Prats, A.C., Vagner, S., 2003. A single internal ribosome entry site containing a G quartet RNA structure drives fibroblast growth factor 2 gene expression at four alternative translation initiation codons. *J. Biol. Chem.* 278, 39330–39336.
- Bonnet, E., Wuyts, J., Rouzé, P., Van de Peer, Y., 2004a. Detection of 91 potential conserved plant microRNAs in *Arabidopsis thaliana* and *Oryza sativa* identifies important target genes. *Proc. Natl. Acad. Sci. USA* 101, 11511–11516.
- Bonnet, E., Wuyts, J., Rouzé, P., Van de Peer, Y., 2004b. Evidence that microRNA precursors, unlike other non-coding RNAs, have lower folding free energies than random sequences. *Bioinformatics* 20, 2911–2917.
- Brantl, S., 2004. Bacterial gene regulation: from transcription attenuation to riboswitches and ribozymes. *Trends Microbiol.* 12, 473–475.
- Breaker, R.R., 2002. Engineered allosteric ribozymes as biosensors components. *Curr. Opin. Biotechnol.* 13, 31–39.
- Brennicke, A., Marchfelder, A., Binder, S., 1999. RNA editing. *FEMS Microbiol. Rev.* 23, 297–316.
- Brown, J., 1999. The ribonuclease P database. *Nucleic Acids Res.* 27, 314.
- Brown, C.J., Ballabio, A., Rupert, J.L., Lafrenière, R.G., Grompe, M., Tonlorenzi, R., Willard, H.F., 1991. A gene from the region of the human X inactivation centre is expressed exclusively from the inactive X chromosome. *Nature* 349, 38–44.
- Buratti, E., Baralle, F.E., 2004. Influence of RNA secondary structure on the pre-mRNA splicing process. *Mol. Cell Biol.* 24, 10505–10514.
- Caetano-Anollés, G., 2002a. Evolved RNA secondary structure and the rooting of the universal tree. *J. Mol. Evol.* 54, 333–345.
- Caetano-Anollés, G., 2002b. Tracing the evolution of RNA structure in ribosomes. *Nucleic Acids Res.* 30, 2575–2587.
- Cai, X., Hagedorn, C.H., Cullen, B.R., 2004. Human microRNAs are processed from capped, polyadenylated transcripts that can also function as mRNAs. *RNA* 10, 1957–1966.
- Candelon, B., Guilloux, K., Ehrlich, S.D., Sorokin, A., 2004. Two distinct groups of rRNA operons in the *Bacillus cereus* group. *Microbiology* 150, 601–611.
- Carmell, M.A., Hannon, G.J., 2004. RNase III enzymes and the initiation of gene silencing. *Nat. Struct. Mol. Biol.* 11, 214–218.
- Caranza, S., Giribet, G., Ribera, C., Baguñà, J., Riutort, M., 1996. Evidence that two types of 18S rDNA coexist in the genome of *Dugesia (Schmidtea) mediterranea* (platyhelminthes, turbellaria, tricladida). *Mol. Biol. Evol.* 13, 824–832.
- Caranza, S., Baguñà, J., Riutort, M., 1999. Origin and evolution of paralogous rRNA gene clusters within the flatworm family dugesiidae (platyhelminthes, tricladida). *J. Mol. Evol.* 49, 250–259.
- Carter, R.J., Dubchak, I., Holbrook, S.R., 2001. A computational approach to identify genes for functional RNAs in genomic sequences. *Nucleic Acids Res.* 29, 3928–3938.
- Cavaillé, J., Buiting, K., Kieffmann, M., Lalande, M., Brennan, C.I., Horsthemke, B., Bachellerie, J.-P., Hüttenthaler, A., 2000. Identification of brain-specific and imprinted small nucleolar RNA genes exhibiting an unusual genomic organization. *Proc. Natl. Acad. Sci. USA* 97, 14311–14316.
- Cavalier-Smith, T., Chao, E.E.-Y., 2003. Phylogeny of Choanozoa, Apusozoa, and other protzoa and the early eukaryote megaevolution. *J. Mol. Evol.* 56, 540–563.
- Cawley, S., Bekiranov, S., Ng, H.H., Kapranov, P., Sekinger, E.A., Kampa, D., Piccolboni, A., Sementchenko, V., Cheng, J., Williams, A.J., Wheeler, R., Wong, B., Drenkow, J., Yamanaka, M., Patel, S., Brubaker, S., Tammana, H., Helt, G., Struhl, K., Gingeras, T.R., 2004. Unbiased mapping of transcription factor binding sites along human chromosomes 21 and 22 points to widespread regulation of noncoding RNAs. *Cell* 116, 499–509.
- Cecconi, F., Crosio, C., Mariottini, P., Cesareni, G., Giorgi, M., Brenner, S., Amaldi, F., 1996. A functional role for some fugu introns larger than the typical short ones: the example of the gene coding for ribosomal protein S7 and snoRNA U17. *Nucleic Acids Res.* 24, 3167–3172.
- Cervelli, M., Cecconi, F., Giorgi, M., Annesi, F., Oliverio, M., Mariottini, P., 2002. Comparative structure analysis of vertebrate U17 small nucleolar RNA (snoRNA). *J. Mol. Evol.* 54, 166–179.
- Cervelli, M., Oliverio, M., Bellini, A., Bologna, M., Cecconi, F., Mariottini, P., 2003. Structural and sequence evolution of U17 small nucleolar RNA (snoRNA) and its phylogenetic congruence in chelonians. *J. Mol. Evol.* 57, 73–84.
- Chartrand, P., Meng, X.H., Singer, R.H., Long, R.M., 1999. Structural elements required for the localization of ASH1 mRNA and of a green fluorescent protein reporter particle *in vivo*. *Curr. Biol.* 9, 333–336.

- Chen, J.-L., Greider, C.W., 2004. An emerging consensus for telomerase RNA structure. *Proc. Natl. Acad. Sci. USA* 101, 14683–14684.
- Chen, J.H., Le, S.Y., Shapiro, B., Currey, K.M., Maizel Jr., J.V., 1990. A computational procedure for assessing the significance of RNA secondary structure. *Comput. Appl. Biosci.* 6, 7–18.
- Chen, J.L., Blasco, M.A., Greider, C.W., 2000. Secondary structure of vertebrate telomerase RNA. *Cell* 100, 503–514.
- Chen, S., Lesnik, E.A., Hall, T.A., Sampath, R., Griffey, R.H., Eker, D., Blyn, L., 2002. A bioinformatics based approach to discover small RNA genes in the *Escherichia coli* genome. *Biosystems* 65, 157–177.
- Chureau, C., Prissette, M., Bourdet, A., Barbe, V., Cattolico, L., Jones, L., Eggen, A., Avner, P., Duret, L., 2002. Comparative sequence analysis of the X-inactivation center region in mouse, human, and bovine. *Genome Res.* 12, 894–908.
- Clayton, C.E., 2002. Life without transcriptional control? From fly to man and back again. *EMBO J.* 21, 1881–1888.
- Clouet d'Orval, B., Bortolin, M.L., Gaspin, C., Bachellerie, J.P., 2001. Box C/D RNA guides for the ribose methylation of archaeal tRNAs. The tRNATrp intron guides the formation of two ribose-methylated nucleosides in the mature tRNATrp. *Nucleic Acids Res.* 29, 4518–4529.
- Collins, L.J., 2004. Lost in the RNA world. Ph.D. Thesis, Allan Wilson Center, Massey University, Palmerston North, New Zealand.
- Collins, L.J., Moulton, V., Penny, D., 2000. Use of RNA secondary structure for studying the evolution of RNase P and RNase MRP. *J. Mol. Evol.* 51, 194–204.
- Collins, L.J., Macke, T.J., Penny, D., 2004. Searching for ncRNAs in eukaryotic genomes: maximizing biological input with RNAmotif. *J. Integrated Bioinform.* 6, 15 <http://journal.imbio.de/>.
- Coventry, A., Kleitman, D.J., Berger, B., 2004. MSARI: multiple sequence alignments for statistical detection of RNA secondary structure. *Proc. Natl. Acad. Sci. USA* 101, 12102–12107.
- Crosthwaite, S.K., 2004. Circadian clocks and natural antisense RNA. *FEBS Lett.* 567, 49–54.
- Crucis, S., Chatterjee, S., Gavis, E.R., 2000. Overlapping but distinct RNA elements control repression and activation of nanos translation. *J. Mol. Cell Biol.* 3, 457–467.
- Dahlberg, J.E., Lund, E., 1988. The genes and transcription of the major small nuclear RNAs. In: Birnstiel, M.L. (Ed.), *Structure and Function of Major and Minor Small Nuclear Ribonucleoprotein Particles*. Springer, Berlin, pp. 38–70.
- Dalphin, E., Stockwell, P.A., Tate, W.P., Brown, C.M., 1999. TransTerm, the translational signal data base, extended to include full coding sequences and untranslated regions. *Nucleic Acids Res.* 27, 293–294.
- Dandjinou, A.T., Lévesque, N., Larose, S., Lucier, J.-F., Elela, S.A., Wellinger, R.J., 2004. A phylogenetically based secondary structure for the yeast telomerase RNA. *Curr. Biol.* 14, 1148–1158.
- Day, W.H.E., Edelsbrunner, H., 1984. Efficient algorithms for agglomerative hierarchical clustering methods. *J. Classifications* 1, 7–24.
- Day, D.A., Tuite, M.F., 1998. Post-transcriptional gene regulatory mechanisms in eukaryotes: an overview. *J. Endocrinol.* 157, 361–371.
- Dayton, E.T., Konings, D.A., Powell, D.M., Shapiro, B.A., Butini, L., Maizel, J.V., Dayton, A.I., 1992. Extensive sequence-specific information throughout the CAR/RRE, the target sequence of the human immunodeficiency virus type 1 Rev protein. *J. Virol.* 66, 1139–1151.
- de la Cruz, J., Vioque, A., 2003. A structural and functional study of plastid RNAs homologous to catalytic bacterial RNase P RNA. *Gene* 321, 47–56.
- de Turris, V., Di Leva, G., Calderola, S., Loreni, F., Amaldi, F., Bozzoni, I., 2004. TOP promoter elements control the relative ratio of intron-encoded snoRNA versus spliced mRNA biosynthesis. *J. Mol. Biol.* 344, 383–394.
- Delihas, N., 2003. Annotation and evolutionary relationships of a small regulatory RNA gene *micF* and its target *ompF* in *Yersinia* species. *BMC Microbiol.* 3, 13 (15 pp.).
- Dennis, P.P., Omer, A., Lowe, T., 2001. A guided tour: small RNA function in archaea. *Mol. Microbiol.* 40, 509–519.
- di Bernardo, D., Down, T., Hubbard, T., 2003. ddbRNA: detection of conserved secondary structures in multiple alignments. *Bioinformatics* 19, 1606–1611.
- Di Giulio, M., 2004. The origin of the tRNA molecule: implications for the origin of protein synthesis. *J. Theor. Biol.* 226, 89–93.
- Dieterich, C., Grossmann, S., Tanzer, A., Ropcke, S., Arndt, P. F., Stadler, P. F., Vingron, M., 2005. Comparative promoter region analysis powered by CORG. *BMC Genom.* in press.
- Dominski, Z., Yang, X.-C., Purdy, M., Marzluff, W.F., 2003. Cloning and characterization of the *Drosophila* U7 small nuclear RNA. *Proc. Natl. Acad. Sci. USA* 100, 9422–9427.

- Domitrovich, A.M., Kunkel, G.R., 2003. Multiple, dispersed human U6 small nuclear RNA genes with varied transcriptional efficiencies. *Nucleic Acids Res.* 31, 2344–2352.
- Doolittle, W.F., Brown, J.R., 1994. Tempo, mode, the progenote, and the universal root. *Proc. Natl. Acad. Sci. USA* 91, 6721–6728.
- Doudna, J.A., Cech, T.R., 2002. The chemical repertoire of natural ribozymes. *Nature* 418, 222–228.
- Duret, L., Dorkeld, F., Gautier, C., 1993. Strong conservation of non-coding sequences during vertebrates evolution: potential involvement in post-transcriptional regulation of gene expression. *Nucleic Acids Res.* 21, 2315–2322.
- Eddy, S.R., 2001. Non-coding RNA genes and the modern RNA world. *Nat. Genet.* 2, 919–929.
- Eddy, S.R., 2002. A memory-efficient dynamic programming algorithm for optimal alignment of a sequence to an RNA secondary structure. *BMC Bioinform.* 3, 18.
- Edvardsson, S., Gardner, P.P., Poole, A.M., Hendy, M.D., Penny, D., Moulton, V., 2003. A search for H/ACA snRNAs in yeast using MFE secondary structure prediction. *Bioinformatics* 19, 865–873.
- Eigen, M., Winkler-Oswatitsch, R., 1981. Transfer-RNA, an early gene? *Naturwissenschaften* 68, 282–292.
- Eigen, M., Lindemann, B.F., Tietze, M., Winkler-Oswatitsch, R., Dress, A.W.M., von Haeseler, A., 1989. How old is the genetic code? Statistical geometry of tRNA provides an answer. *Science* 244, 673–679.
- Elbashir, W., Lendeckel, S., Tuschl, T., 2001. RNA interference is mediated by 21- and 22-nucleotide RNAs. *Genes Dev.* 15, 188–200.
- Enright, C.A., Maxwell, E.S., Elicieri, G.L., Sollner-Webb, B., 1996. 5'ETS rRNA processing facilitated by four small RNAs: U14, E3, U17, and U3. *RNA* 2, 1094–1099.
- Enright, A.J., John, B., Gaul, U., Tuschl, T., Sander, C., Marks, D.S., 2003. MicroRNA targets in *Drosophila*. *Genome Biol.* 5 (1) (Article R1).
- Erdmann, V., Szymański, M., Hochberg, A., de Groot, N., Barciszewski, J., 1999. Collection of mRNA-like non-coding RNAs. *Nucleic Acids Res.* 27, 192–195.
- Erdmann, V.A., Szymański, M., Hochberg, A., deGroot, N., Barciszewski, J., 2000. Non-coding, mRNA-like RNAs database Y2K. *Nucleic Acids Res.* 28, 197–2000.
- Erdmann, V., Barciszewska, M., Hochberg, A., de Groot, N., Barciszewski, J., 2001. Regulatory RNAs. *Cell. Mol. Life Sci.* 58, 960–977.
- Escriva, H., Manzon, L., Youson, J., Laudet, V., 2002. Analysis of lamprey and hagfish genes reveals a complex history of gene duplications during early vertebrate evolution. *Mol. Biol. Evol.* 19, 1440–1450.
- Estévez, A.M., Simpson, L., 1999. Uridine insertion/deletion RNA editing in trypanosome mitochondria—a review. *Gene* 240, 247–260.
- Fabre, E., Dujon, B., Richard, G., 2002. Transcription and nuclear transport of CAG/CTG trinucleotide repeats in yeast. *Nucleic Acids Res.* 30, 3540–3547.
- Fagegaltier, D., Lescure, A., Walczak, R., Carbon, P., Krol, A., 2000. Structural analysis of new local features in SECIS RNA hairpins. *Nucleic Acids Res.* 28 (14), 2679–2689.
- Farris, A.D., Koelsch, G., Pruijn, G.J., van Venrooij, W.J., Harley, J.B., 1999. Conserved features of Y RNAs revealed by automated phylogenetic secondary structure analysis. *Nucleic Acids Res.* 27, 1070–1078.
- Fayat, G., Mayaux, F.J., Sacerdot, C., Fromant, M., Springer, M., Grunberg-Manago, M., Blanquet, S., 1983. *Escherichia coli* phenylalanyl-tRNA synthetase operon region. Evidence for an attenuation mechanism. Identification of the gene for the ribosomal protein L20. *J. Mol. Biol.* 171, 239–261.
- Felden, B., Massire, C., Westhof, E., Atkins, J.F., Gesteland, R.F., 2001. Phylogenetic analysis of tmRNA genes within a bacterial subgroup reveals a specific structural signature. *Nucleic Acids Res.* 29, 1602–1607.
- Ferreira, M.G., Miller, K.M., Cooper, J.P., 2004. Indecent exposure: when telomeres become uncapped. *Mol. Cell* 13, 7–18.
- Filippov, V., Solovyev, V., Filippova, M., Gill, S., 2000. A novel type of RNase III family proteins in eukaryotes. *Gene* 245, 213–221.
- Flamm, C., Hofacker, I.L., Maurer-Stroh, S., Stadler, P.F., Zehl, M., 2000. Design of multi-stable RNA molecules. *RNA* 7, 254–265.
- Franke, A., Baker, B., 2000. Dosage compensation rox!. *Curr. Opin. Cell Biol.* 12, 351–354.
- Freeland, S.J., Knight, R.D., Landweber, L.F., 1999. Do proteins predate DNA? *Science* 286, 690–692.
- Frenkel, F.E., Chaley, M.B., Korotkov, E.V., Skryabin, K.G., 2004. Evolution of tRNA-like sequences and genome variability. *Gene* 335, 57–71.
- Gardner, P.P., Giegerich, R., 2004. A comprehensive comparison of comparative RNA structure prediction approaches. *BMC Bioinform.* 5, 140.
- Gaudin, C., Zhou, X., Williams, K.P., Felden, B., 2002. Two-piece tmRNA in cyanobacteria and its structural analysis. *Nucleic Acids Res.* 30, 2018–2024.

- Gautheret, D., Major, F., Cedergren, R., 1990. Pattern searching/alignment with RNA primary and secondary structures: an effective descriptor for tRNA. *Comput. Appl. Biosci.* 6, 325–331.
- Gautheret, D., Lambert, A., 2001. Direct RNA motif definition and identification from multiple sequence alignments using secondary structure profiles. *J. Mol. Biol.* 313, 1003–1011.
- Gebauer, F., Hentze, M.W., 2004. Molecular mechanisms of translational control. *Nat. Rev. Mol. Cell Biol.* 5, 827–835.
- Gesteland, R.F., Atkins, J.F. (Eds.), 1993. The RNA World. Cold Spring Harbor Laboratory Press, Plainview, NY.
- Gilbert, W., 1986. The RNA world. *Nature* 319, 618.
- Giles, K.E., Caputi, M., Beemon, K., 2004. Packaging and reverse transcription of snRNAs by retroviruses may generate pseudogenes. *RNA* 10, 299–307.
- Gilley, J., Fried, M., 1998. Evolution of U24 and U36 snoRNAs encoded within introns of vertebrate *rpL7a* gene homologs: unique features of mammalian U36 variants. *DNA Cell Biol.* 17, 591–602.
- Gilmartin, G.M., Schaufele, F., Schaffner, G., Birnstiel, M.L., 1988. Functional analysis of the sea urchin U7 small nuclear RNA. *Mol. Cell Biol.* 8, 1076–1084.
- Gonzalez, I.L., Sylvester, J.E., 2001. Human rDNA: evolutionary patterns within the genes and tandem arrays derived from multiple chromosomes. *Genomics* 73, 255–263.
- Goodfellow, I.G., Kerrigan, D., Evans, D.J., 2003. Structure and functional analysis of the poliovirus *cis*-acting replication element (CRE). *RNA* 9, 124–137.
- Gorodkin, J., Heyer, L.J., Stormo, G.D., 1997. Finding the most significant common sequence and structure motifs in a set of RNA sequences. *Nucleic Acids Res.* 25 (18), 3724–3732.
- Gorodkin, J., Knudsen, B., Zwieb, C., Samuelsson, T., 2001a. SRPDB (signal recognition particle database). *Nucleic Acids Res.* 29, 169–170.
- Gorodkin, J., Stricklin, S.L., Stormo, G.D., 2001b. Discovering common stem-loop motifs in unaligned RNA sequences. *Nucleic Acids Res.* 29 (10), 2135–2144.
- Gott, J.M., Emeson, R.B., 2000. Functions and mechanisms of RNA editing. *Annu. Rev. Genet.* 34, 499–531.
- Gottesman, S., 2004. The small RNA regulators of *Escherichia coli*: roles and mechanisms. *Annu. Rev. Microbiol.* 58, 303–328.
- Gottlob-McHugh, S.G., Levesque, M., MacKenzie, K., Olson, M., Yarosh, O., Johnson, D.A., 1990. Organization of the 5S rRNA genes in the soybean *Glycine max* (L.) Merrill and conservation of the 5S rRNA repeat structure in higher plants. *Genome* 33, 486–494.
- Gräf, S., Strothmann, D., Kurtz, S., Steger, G., 2001. HyPaLib: a database of RNAs and RNA structural elements defined by hybrid patterns. *Nucleic Acids Res.* 29, 196–198.
- Griffiths-Jones, S., 2004. The microRNA registry. *Nucleic Acids Res.* 32 (Database issue), D109–D111.
- Griffiths-Jones, S., Bateman, A., Marshall, M., Khanna, A., Eddy, S., 2003. Rfam: an RNA family database. *Nucleic Acids Res.* 31, 439–441.
- Griffiths-Jones, S., Moxon, S., Marshall, M., Khanna, A., Eddy, S.R., Bateman, A., 2005. Rfam: annotating non-coding RNAs in complete genomes. *Nucleic Acids Res.* 33 (Database issue), 121–124.
- Grummt, I., 2003. Life on a planet of its own: regulation of RNA polymerase i transcription in the nucleolus. *Genes Dev.* 17, 1691–1702.
- Grundy, F.J., Lehman, S.C., Henkin, T.M., 2003. The L box regulon: lysine sensing by leader RNAs of bacterial lysine biosynthesis genes. *Proc. Natl. Acad. Sci. USA* 100, 12057–12062.
- Grüner, W., Giegerich, R., Strothmann, D., Reidys, C., Weber, J., Hofacker, I.L., Stadler, P.F., Schuster, P., 1996a. Analysis of RNA sequence structure maps by exhaustive enumeration, I: neutral networks. *Monatsh. Chem.* 127, 355–374.
- Grüner, W., Giegerich, R., Strothmann, D., Reidys, C., Weber, J., Hofacker, I.L., Stadler, P.F., Schuster, P., 1996b. Analysis of RNA sequence structure maps by exhaustive enumeration, II: structures of neutral networks and shape space covering. *Monatsh. Chem.* 127, 375–389.
- Gueneau de Novoa, P., Williams, K.P., 2004. The tmRNA website: reductive evolution of tmRNA in plastids and other endosymbionts. *Nucleic Acids Res.* 32 (Database issue), D104–D108.
- Guo, P., 2002. Structure and function of phi29 hexameric RNA that drives the viral DNA packaging motor: review. *Prog. Nucleic Acid Res. Mol. Biol.* 72, 415–472.
- Gürsoy, H.-C., Koper, D., Benecke, B.-J., 2000. The vertebrate 7S K RNA separates hagfish (*Myxine glutinosa*) and lamprey (*Lampetra fluviatilis*). *J. Mol. Evol.* 50, 456–464.
- Gustafson, A.M., Allen, E., Givan, S., Smith, D., Carrington, J.C., Kasschau, K.D., 2005. ASRP: the *Arabidopsis* Small RNA Project Database. *Nucleic Acids Res.* 33, D637–D640.
- Haas, E.S., Banta, A.B., Harris, J.K., Pace, N.R.P., Brown, J.W., 1996. Structure and evolution of ribonuclease P RNA in Gram-positive bacteria. *Nucleic Acids Res.* 24, 4775–4782.

- Hackermüller, J., Meisner, N.-C., Auer, M., Jaritz, M., Stadler, P.F., 2005. The effect of RNA secondary structures on RNA-ligand binding and the modifier RNA mechanism: a quantitative model. *Gene* 345, 3–12.
- Haebel, P.W., Gutmann, S., Ban, N., 2004. Dial tm for rescue: tmRNA engages ribosomes stalled on defective mRNAs. *Curr. Opin. Struct. Biol.* 14, 58–65.
- Hannon, G.J., 2002. RNA interference. *Nature* 418, 244–251.
- Harris, R.J., Elder, D., 2000. Ribozyme relationships: the hammerhead, hepatitis delta, and hairpin ribozymes have a common origin. *J. Mol. Evol.* 51, 182–184.
- Hartmann, E., Hartmann, R.K., 2003. The enigma of ribonuclease P evolution. *Trends Genet.* 19, 561–569.
- Havgaard, J.H., Lingsø, R., Stormo, G.D., Gorodkin, J., 2005. Pairwise local structural alignment of RNA sequences with sequence similarity less than 40%. *Bioinformatics (Epub Januray 18 2005)*.
- Henkin, T.M., Yanofsky, C., 2002. Regulation by transcription attenuation in bacteria: how RNA provides instructions for transcription termination/antitermination decision. *BioEssays* 24, 700–707.
- Henras, A.K., Dez, C., Henry, Y., 2004. RNA structure and function in C/D and H/ACA s(no)RNAs. *Curr. Opin. Struct. Biol.* 14, 335–343.
- Hentze, M.W., Kühn, L.C., 1996. Molecular control of vertebrate iron metabolism: mRNA-based regulatory circuits operated by iron, nitric oxide, and oxidative stress. *Proc. Natl. Acad. Sci. USA* 93, 8175–8182.
- Hernandez, N., 2001. Small nuclear RNA genes: a model system to study fundamental mechanisms of transcription. *J. Biol. Chem.* 276, 26733–26736.
- Hershberg, R., Altuvia, S., Margalit, H., 2003. A survey of small RNA-encoding genes in *Escherichia coli*. *Nucleic Acids Res.* 31, 1813–1820.
- Hesketh, J., 2004. 3'-untranslated regions are important in mRNA localization and translation: lessons from selenium and metallothionein. *Biochem. Soc. Trans.* 32, 990–993.
- Higgs, P.G., Jameson, D., Jow, H., Rattray, M., 2003. The evolution of tRNA-leu genes in animal mitochondrial genomes. *J. Mol. Evol.* 435–445.
- Hillis, D.M., Dixon, M.T., 1991. Ribosomal DNA: molecular evolution and phylogenetic inference. *Q. Rev. Biol.* 66, 411–453.
- Hinz, S., Göringer, H.U., 1999. The guide RNA database (3.0). *Nucleic Acids Res.* 27, 168.
- Hobert, O., 2004. Common logic of transcription factor and microRNA action. *Trends Biochem. Sci.* 29, 462–468.
- Höchsmann, M., Töller, T., Giegerich, R., Kurtz, S., 2003. Local similarity in RNA secondary structures. In: Proceedings of the Computational Systems Bioinformatics Conference, Stanford, CA, August 2003 (CSB 2003), pp. 159–168.
- Hofacker, I.L., 2003. Vienna RNA secondary structure server. *Nucleic Acids Res.* 31, 3429–3431.
- Hofacker, I.L., Stadler, P.F., 1999. Automatic detection of conserved base pairing patterns in RNA virus genomes. *Comput. Chem.* 23, 401–414.
- Hofacker, I.L., Fontana, W., Stadler, P.F., Bonhoeffer, L.S., Tacker, M., Schuster, P., 1994. Fast folding and comparison of RNA secondary structures. *Monatsh. Chem.* 125, 167–188.
- Hofacker, I.L., Fekete, M., Flamm, C., Huynen, M.A., Rauscher, S., Stolorz, P.E., Stadler, P.F., 1998. Automatic detection of conserved RNA structure elements in complete RNA virus genomes. *Nucleic Acids Res.* 26, 3825–3836.
- Hofacker, I.L., Fekete, M., Stadler, P.F., 2002. Secondary structure prediction for aligned RNA sequences. *J. Mol. Biol.* 319, 1059–1066.
- Hofacker, I.L., Bernhart, S.H.F., Stadler, P.F., 2004a. Alignment of RNA base pairing probability matrices. *Bioinformatics* 20, 2222–2227.
- Hofacker, I.L., Stocsits, R., Stadler, P.F., 2004b. Conserved RNA secondary structures in viral genomes: a survey. *Bioinformatics* 20, 1495–1499.
- Holland, P.W.H., García-Fernández, J., Williams, N.A., Sidow, A., 1994. Gene duplication and the origins of vertebrate development. *Development (Suppl.)*, 125–133.
- Holmes, I., 2004. A probabilistic model for the evolution of RNA structure. *BMC Bioinformatics* 5, 166.
- Hong, M., Simpson, L., 2003. Genomic organization of *Trypanosoma brucei* kinetoplast DNA minicircles. *Prostist* 154, 265–279.
- Hopper, A.K., Phizicky, E.M., 2003. tRNA transfers to the limelight. *Genes Dev.* 17, 162–180.
- Hu, Y.-J., 2002. Prediction of consensus structural motifs in a family of coregulated RNA sequences. *Nucleic Acids Res.* 30, 3886–3893.
- Hu, Y., 2003. GPRM: a genetic programming approach to finding common RNA secondary structure elements. *Nucleic Acids Res.* 31, 3446–3449.

- Huang, Z.P., Zhou, H., Liang, D., Qu, L.H., 2004. Different expression strategy: multiple intronic gene clusters of box H/ACA snoRNA in *Drosophila melanogaster*. *J. Mol. Biol.* 341, 669–683.
- Hudelot, C., Gowri-Shankar, V., Jow, H., Rattray, M., Higgs, P.G., 2003. RNA-based phylogenetic methods: application to mammalian mitochondrial RNA sequences. *Mol. Phylogenet. Evol.* 28, 241–252.
- Huez, I., Créancier, L., Audigier, S., Gensac, M., Prats, A., Prats, H., 1998. Two independent internal ribosome entry sites are involved in translation initiation of vascular endothelial growth factor mRNA. *Mol. Cell. Biol.* 18, 6178–6190.
- Hüttenhofer, A., Kiefmann, M., Neier-Ewert, S., O'Brien, J., Lehrach, H., Bachellerie, J., Brosius, J., 2001. Rnomics: an experimental approach that identifies 201 candidates for novel, small, non-messenger RNAs in mouse. *EMBO J.* 20, 2943–2953.
- Huynen, M.A., Stadler, P.F., Fontana, W., 1996. Smoothness within ruggedness: the role of neutrality in adaptation. *Proc. Natl. Acad. Sci. USA* 93, 397–401.
- Illangasekare, M., Yarus, M., 1999. A tiny RNA that catalyzes both aminoacyl-RNA and peptidyl-RNA synthesis. *RNA* 5, 1482–1489.
- Iyer, V.R., 2004. Exploring the post-transcriptional RNA world with DNA microarrays. *Trends Biotechnol.* 22, 498–500.
- Jacob, Y., Seif, E., Paquet, P.-O., Lang, F.B., 2004. Loss of the mRNA-like region in mitochondrial tmRNAs of jakobids. *RNA* 10, 605–614.
- Jadhav, V.R., Yarus, M., 2002. Coenzymes as coribozymes. *Biochimie* 84, 877–888.
- Jády, B.E., Kiss, T., 2001. A small nucleolar guide RNA functions both in 2'-O-methylation and pseudouridylation of U5 spliceosomal RNA. *EMBO J.* 20, 541–551.
- Jády, B.E., Bertrand, E., Kiss, T., 2004. Human telomerase RNA and box H/ACA scaRNAs share a common Cajal body specific localization signal. *J. Cell Biol.* 164, 647–652.
- Jameson, D., Gibson, A.P., Hudelot, C., Higgs, P.G., 2003. OGRe: a relational database for comparative analysis of mitochondrial genomes. *Nucleic Acids Res.* 31, 202–206.
- Jang, S.K., Krausslich, H.G., Nicklin, M.J., Duke, G.M., Palmenberg, A.C., Wimmer, E., 1988. A segment of the 5' nontranslated region of encephalomyocarditis virus RNA directs internal entry of ribosomes during *in vitro* translation. *J. Virol.* 62, 2636–2643.
- Jareborg, N., Birney, E., Durbin, R., 1999. Comparative analysis of noncoding regions of 77 orthologous mouse and human gene pairs. *Genome Res.* 9, 815–824.
- Jeffares, D., Poole, A.M., Penny, D., 1998. Relics from the RNA world. *J. Mol. Evol.* 46, 18–36.
- Ji, Y., Xing, X., Stormo, G.D., 2004. A graph theoretical approach for predicting common RNA secondary structure motifs including pseudoknots in unaligned sequences. *Bioinformatics* 20, 1591–1602.
- Johnston, W.K., Unrau, P.J., Lawrence, M.J., Glasner, M.E., Bartel, D.P., 2001. RNA-catalyzed RNA polymerization: accurate and general RNA-templated primer extension. *Science* 292, 1319–1325.
- Jones-Roades, M.W., Bartel, D.P., 2004. Computational identification of plant microRNAs and their targets, including a stress-induced miRNA. *Mol. Cell* 14, 787–799.
- Jow, H., Hudelot, C., Rattray, M., Higgs, P.G., 2002. Bayesian phylogenetics using an RNA substitution model applied to early mammalian evolution. *Mol. Biol. Evol.* 19, 1591–1601.
- Joyce, G.F., 2002. The antiquity of RNA-based evolution. *Nature* 418, 214–221.
- Joyce, G.F., 2004. Directed evolution of nucleic acid enzymes. *Annu. Rev. Biochem.* 73, 791–836.
- Juan, V., Crain, C., Wilson, C., 2000. Evidence for evolutionarily conserved secondary structure in the H19 tumour suppressor RNA. *Nucleic Acids Res.* 28, 1221–1227.
- Kampa, D., Cheng, J., Kapranov, P., Yamanaka, M., Brubaker, S., Cawley, S., Drenkow, J., Piccolboni, A., Bekiranov, S., Helt, G., Tammana, H., Gingras, T.R., 2004. Novel RNAs identified from an in-depth analysis of the transcriptome of human chromosomes 21 and 22. *Genome Res.* 14, 331–342.
- Katz, L., Burge, C.B., 2003. Widespread selection for local RNA secondary structure in coding regions of bacterial genes. *Genome Res.* 13, 2042–2051.
- Keenan, R.J., Freyman, D.M., Stroud, R.M., Walter, P., 2001. The signal recognition particle. *Annu. Rev. Biochem.* 70, 755–775.
- Keiler, K.C., Shapiro, L., Williams, K.P., 2000. tmRNAs that encode proteolysis-inducing tags are found in all known bacterial genomes: a two-piece tmRNA functions in *caulobacter*. *Proc. Natl. Acad. Sci. USA* 97, 7778–7783.
- Kelleher, C., Teixeira, M.T., Förstemann, K., Lingner, J., 2002. Telomerase: biochemical considerations for enzyme and substrate. *Trends Biochem. Sci.* 27, 572–579.

- Khaitovich, P., Mankin, A.S., Green, R., Lancaster, L., Noller, H.F., 1999. Characterization of functionally active subribosomal particles from *Thermus aquaticus*. Proc. Natl. Acad. Sci. USA 96, 85–90.
- Kidner, C.A., Martienssen, R.A., 2005. The developmental role of microRNA in plants. Curr. Opin. Plant Biol. 8, 38–44.
- Kiss, T., 2001. Small nucleolar RNA-guided post-transcriptional modification of cellular RNAs. EMBO J. 20, 3617–3622.
- Klein, R.J., Eddy, S.R., 2003. RESEARCH: finding homologs of single structured RNA sequences. BMC Bioinform. 4 (44), 1471–2105.
- Klein, R.J., Misulovin, Z., Eddy, S.R., 2002. Noncoding RNA genes identified in AT-rich hyperthermophiles. Proc. Natl. Acad. Sci. USA 99, 7542–7547.
- Knudsen, B., Hein, J.J., 1999. Using stochastic context free grammars and molecular evolution to predict RNA secondary structure. Bioinformatics 15, 446–454.
- Knudsen, B., Hein, J., 2003. Pfold: RNA secondary structure prediction using stochastic context-free grammars. Nucleic Acids Res. 31, 3423–3428.
- Kohtz, J., Fishell, G., 2004. Developmental regulation of EVF-1 a novel non-coding RNA transcribed upstream of the mouse Dlx6 gene. Gene Exp. Patterns 4, 407–412.
- Komatsu, Y., 2004. Regulation of ribozyme activity with short oligonucleotides. Biol. Pharma. Bull. 27, 457–462.
- Konecny, J., Schöninger, M., Hofacker, I.L., Weitze, M.-D., Hofacker, G.L., 2000. Concurrent neutral evolution of mRNA secondary structures and encoded proteins. J. Mol. Evol. 50, 238–242.
- Koper-Emde, D., 2004. Phylogenetische Heterogenität der 7S-RNAs von Eukaryonten. Ph.D. Thesis, University of Bochum.
- Korecnic, D., Ahel, I., Schelert, J., Sacher, M., Ruan, B., Stathopoulos, C., Blum, P., Ibba, M., Söll, D., 2004. A freestanding proofreading domain is required for protein synthesis quality control in archaea. Proc. Natl. Acad. Sci. USA 101, 10260–10265.
- Kozak, M., 1989. The scanning model for translation: an update. J. Cell. Biol. 108, 229–241.
- Kozak, M., 1991. An analysis of vertebrate mRNA sequences: intimations of translational control. J. Cell Biol. 115 (4), 887–903.
- Kozak, M., 2001. New ways of initiating translation in eukaryotes? Mol. Cell. Biol. 21 (6), 1899–1907.
- Krasilnikov, A.S., Xiao, Y., Pan, T., Mondragón, A., 2004. Basis for structural diversity in homologous RNAs. Science 306, 104–107.
- Krol, A., 2002. Evolutionarily different RNA motifs and RNA–protein complexes to achieve selenoprotein synthesis. Biochimie 84, 765–774.
- Kwek, K.Y., Murphy, S., Furger, A., Thomas, B., O’Gorman, W., Kimura, H., Proudfoot, N.J., Akoulitchev, A., 2002. U1 snRNA associates with TFIIH and regulates transcriptional initiation. Nat. Struct. Biol. 9, 800–805.
- Lafontaine, D., Tollervey, D., 2002. Birth of the snoRNPs: the evolution of the modification-guide snoRNAs. Trends Biochem. Sci. 23, 383–388.
- Laforest, M.-J., Bullerwell, C.E., Forget, L., Lang, F.B., 2004. Origin, evolution, and mechanism of 5’tRNA editing in chytridiomycete fungi. RNA 10, 1191–1199.
- Lagos-Quintana, M., Rauhut, R., Lendeckel, W., Tuschl, T., 2001. Identification of novel genes coding for small expressed RNAs. Science 294, 853–857.
- Lagos-Quintana, M., Rauhut, R., Meyer, J., Borkhardt, A., Tuschl, T., 2003. New microRNAs from mouse and human. RNA 9, 175–179.
- Lai, E.C., Tomancak, P., Williams, R.W., Rubin, G.M., 2003. Computational identification of *Drosophila* microRNA genes. Genome Biol. 4, R42.
- Landry, J., Medstrand, P., Mager, D.L., 2001. Repetitive elements in the 5’ untranslated region of a human zinc-finger gene modulate transcription and translation efficiency. Genomics 76 (1–3).
- Landweber, L.F., 1992. The evolution of RNA editing in kinetoplastid protozoa. Biosystems 28, 41–45.
- Landweber, L.F., Gilbert, W., 1994. Phylogenetic analysis of RNA editing: a primitive genetic phenomenon. Proc. Natl. Acad. Sci. USA 91, 918–921.
- Lang, B.F., Gray, M.W., Burger, G., 1999. Mitochondrial genome evolution and the origin of eukaryotes. Annu. Rev. Genet. 33, 351–397.
- Lariza, A., Makalowski, W., Pesole, G., Saccone, G., 2002. Evolutionary dynamics of mammalian mRNA untranslated regions by comparative analysis of orthologous human, artiodactyl and rodent gene pairs. Comput. Chem. 26, 479–490.
- Laslett, D., Canback, B., Andersson, S., 2002. BRUCE: a program for the detection of transfer-messenger RNA genes in nucleotide sequences. Nucleic Acids Res. 30, 3449–3453.

- Lau, N.C., Lim, L.P., Weinstein, E.G., Bartel, D.P., 2001. An abundant class of tiny RNAs with probable regulatory roles in *Caenorhabditis elegans*. *Science* 294, 858–862.
- Lavorgna, G., Dahary, D., Lehner, B., Sorek, R., Sanderson, C.M., Casari, G., 2004. In search of antisense. *Trends Biochem. Sci.* 29.
- Le, S., Maizel, J.V., 1997. A common RNA structural motif involved in the internal initiation of translation of cellular mRNAs. *Nucleic Acids Res.* 25, 362–369.
- Le, S.V., Chen, J.H., Currey, K.M., Maizel Jr., J.V., 1988. A program for predicting significant RNA secondary structures. *Comput. Appl. Biosci.* 4, 153–159.
- Le, S.Y., Zhang, K., Maizel Jr., J.V., 2002. RNA molecules with structure dependent functions are uniquely folded. *Nucleic Acids Res.* 30, 3574–3582.
- Le, S.Y., Chen, J.H., Konings, D., Maizel Jr., J.V., 2003. Discovering well-ordered folding patterns in nucleotide sequences. *Bioinformatics* 19, 354–361.
- LeCuyer, K.A., Crothers, D.M., 1994. Kinetics of an RNA conformational switch. *Proc. Natl. Acad. Sci. USA* 91, 3373–3377.
- Lee, R., Ambros, V., 2001. An extensive class of small RNAs in *Caenorhabditis elegans*. *Science* 294, 862–864.
- Lee, J.T., Davidow, L.S., Warshawsky, D., 1999. *Tsix*, a gene antisense to *Xist* at the X-inactivation centre. *Nat. Genet.* 21, 400–404.
- Lee, N., Bessho, Y., Wei, K., Szostak, J.W., Suga, H., 2000. Ribozyme-catalyzed tRNA aminoacylation. *Nat. Struct. Biol.* 7, 28–33.
- Lee, Y., Jeon, K., Lee, J.T., Kim, S., Kim, V.N., 2002. MicroRNA maturation: stepwise processing and subcellular localization. *EMBO J.* 21, 4663–4670.
- Lee, Y., Ahn, C., Han, J., Choi, H., Kim, J., Yim, J., Lee, J., Provost, P., Rådmark, O., Kim, S., Kim, V.N., 2003. The nuclear RNase III Drosha initiates microRNA processing. *Nature* 425, 415–419.
- Lee, Y., Kim, M., Han, J., Yeom, K.H., Lee, S., Baek, S.H., Kim, V.K., 2004. MicroRNA genes are transcribed by RNA polymerase II. *EMBO J.* 23, 4051–4060.
- Legendre, M., Lambert, A., Gautheret, D., 2005. Profile-based detection of microRNA precursors in animal genomes. *Bioinformatics* (Epub ahead of print).
- Lesnik, E.A., Fogel, G.B., Weekes, D., Henderson, T.J., Levene, H.B., Sampath, R., Ecker, D.J., 2005. Identification of conserved regulatory RNA structures in prokaryotic metabolic pathway genes. *Biosystems* (doi:10.1016/j.biosystems.2004.11.002).
- Li, K., Williams, R.S., 1995. Cloning and characterization of three new murine genes encoding short homologues of RNase P RNA. *J. Biol. Chem.* 270, 25281–25285.
- Li, Y., Altman, S., 2004. In search of RNase P RNA from microbial genomes. *RNA* 10, 1533–1540.
- Liang, X.H., Xu, Y.X., Michaeli, S., 2002. The spliced-leader associated RNA is a trypanosome-specific sn(o)RNA that has the potential to guide pseudouridine formation on SL RNA. *RNA* 8, 237–246.
- Liao, D., 1999. Concerted evolution: molecular mechanisms and biological implications. *Am. J. Hum. Genet.* 64, 24–30.
- Liao, D., Pavelitz, T., Kidd, J.R., Kidd, K.K., Weiner, A.M., 1997. Concerted evolution of the tandemly repeated genes encoding human U2 snRNA (the RNU2 locus) involves rapid intrachromosomal homogenization and rare interchromosomal gene conversion. *EMBO J.* 16, 588–598.
- Lilley, D.M.J., 2003. The origins of RNA catalysis in ribozymes. *Trends Biochem. Sci.* 28, 495–501.
- Lin, J., Ly, H., Hussain, A., Abraham, M., Pearl, S., Tzfati, Y., Parslow, T.G., Blackburn, E.H., 2004. A universal telomerase RNA core structure includes structured motifs required for binding the telomerase reverse transcriptase protein. *Proc. Natl. Acad. Sci. USA* 101, 14713–14718.
- Lingner, J., Cooper, J.P., Cech, T.R., 1995. Telomerase and DNA end replication: no longer a lagging strand problem? *Science* 269, 1533–1534.
- Lipman, D.J., 1997. Making (anti)sense of non-coding sequence conservation. *Nucleic Acids Res.* 25 (18), 3580–3583.
- Little, R.D., Braaten, B.C., 1989. Genomic organization of human 5 S rDNA and sequence of one tandem repeat. *Genomics* 4, 376–383.
- Liu, C., Bai, B., Skogerbø, G., Cai, L., Deng, W., Zhang, Y., Bu, D., Zhao, Y., Chen, R., 2005. NONCODE: an integrated knowledge database of non-coding RNAs. *Nucleic Acids Res.* 33 (Database issue), D112–D115.
- Lowe, T., Eddy, S., 1997. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* 25, 955–964.
- Lowe, T.M., Eddy, S.R., 1999. A computational screen for methylation guide snoRNAs in yeast. *Science* 19, 1168–1171.

- Lu, S., Cullen, B.R., 2004. Adenovirus VA1 noncoding RNA can inhibit small interfering RNA and MicroRNA biogenesis. *J. Virol.* 78, 12868–12876.
- Luciano, D.J., Mirsky, H., Vendetti, N.J., Maas, S., 2004. RNA editing of a miRNA precursor. *RNA* 10, 1174–1177.
- Lück, R., Steger, G., Riesner, D., 1996. Thermodynamic prediction of conserved secondary structure: application to the RRE element of HIV, the tRNA-like element of CMV, and the mRNA of prion protein. *J. Mol. Biol.* 258, 813–826.
- Lück, R., Gräf, S., Steger, G., 1999. Construct: a tool for thermodynamic controlled prediction of conserved secondary structure. *Nucleic Acids Res.* 27, 4208–4217.
- Lue, N.F., 2004. Adding to the ends: what makes telomerase processive and how important is it? *Bioessays* 26, 955–962.
- Lynch, M., Conery, J.S., 2000. The evolutionary fate and consequences of duplicate genes. *Science* 290, 1151–1155.
- Maas, S., Rich, A., 2000. Changing genetic information through RNA editing. *BioEssays* 22, 790–802.
- MacIntosh, G.C., Wilkerson, C., Green, P.J., 2001. Identification and analysis of *Arabidopsis* expressed sequence tags characteristic of non-coding RNAs. *Plant Physiol.* 127, 765–776.
- Macke, T.J., Ecker, D.J., Gutell, R.R., Gautheret, D., Case, D.A., Sampath, R., 2001. RNAMotif, an RNA secondary structure definition and search algorithm. *Nucleic Acids Res.* 29, 4724–4735.
- Maden, B.E.H., 1990. The numerous modified nucleotides in eukaryotic ribosomal RNA. *Prog. Nucl. Acid Res. Mol. Biol.* 39, 241–303.
- Maidak, B., Cole, J., Lilburn, T., Parker Jr., C., Saxman, P., Farris, R., Garrity, G., Olsen, G., Schmidt, T., Tiedje, J., 2001. The RDP-II (ribosomal database project). *Nucleic Acids Res.* 29, 173–174.
- Malim, M.H., Hauber, J., Le, S.Y., Maizel, J.V., Cullen, B., 1989. The HIV-1 rev trans-activator acts through a structured target sequence to activate nuclear export of unspliced viral mRNA. *Nature* 338, 254–257.
- Mallatt, J.M., Garey, J.R., Shultz, J.W., 2004. Ecdysozoan phylogeny and bayesian inference: first use of nearly complete 28S and 18S rRNA gene sequences to classify the arthropods and their kin. *Mol. Phylogenet. Evol.* 31, 178–191.
- Márquez, L.M., Miller, D.J., MacKenzie, J.B., van Oppen, M.J.H., 2003. Pseudogenes contribute to the extreme diversity of nuclear ribosomal DNA in the hard coral *Acropora*. *Mol. Biol. Evol.* 20, 1077–1086.
- Martineau, Y., Le Bec, C., Monbrun, L., Allo, V., Chiu, I.M., Danos, O., Moine, H., Prats, H., Prats, A.C., 2004. Internal ribosome entry site structural motifs conserved among mammalian fibroblast growth factor 1 alternatively spliced mRNAs. *Mol. Cell. Biol.* 24, 7622–7635.
- Mathews, M.B., 1995. Structure, function, and evolution of adenovirus virus-associated RNAs. *Curr. Top. Microbiol. Immunol.* 199, 173–187.
- Mathews, D.H., Turner, D.H., 2002. Dynalign: an algorithm for finding secondary structures common to two RNA sequences. *J. Mol. Biol.* 317, 191–203.
- Mattick, J.S., 2003. Challenging the dogma: the hidden layer of non-protein-coding RNAs in complex organisms. *Bioessays* 25, 930–939.
- Mattick, J.S., 2004. RNA regulation: a new genetics? *Nat. Genet.* 5, 316–323.
- McCaskill, J.S., 1990. The equilibrium partition function and base pair binding probabilities for RNA secondary structure. *Biopolymers* 29, 1105–1119.
- McCutcheon, J.P., Eddy, S.R., 2003. Computational identification of non-coding RNAs in *Saccharomyces cerevisiae* by comparative genomics. *Nucleic Acids Res.* 31, 4119–4128.
- Meisner, N.-C., Hackermüller, J., Uhl, V., Aszödi, A., Jaritz, M., Auer, M., 2004. mRNA openers and closers: a methodology to modulate AU-rich element controlled mRNA stability by a molecular switch in mRNA conformation. *Chembiochemistry* 5, 1432–1447.
- Mendez, R., Richter, J.D., 2001. Translational control by CPEB: a means to the end. *Nat. Rev. Mol. Cell Biol.* 2, 521–529.
- Merino, E., Yanofsky, C., 2002. Regulation by termination–antitermination: a genomic approach. In: Sonenschein, A.L., Hoch, J.L., Losick, R. (Eds.), *Bacillus subtilis* and its Closest Relatives: from Genes to Cells. ASM Press, Washington DC, 2000, pp. 323–336.
- Michels, A.A., Fraldi, A., Li, Q., Adamson, T.E., Bonnet, F., Nguyen, V.T., Sedore, S.C., Price, J.P., Price, D.H., Lania, L., Bensaude, O., 2004. Binding of the 7SK snRNA turns the HEXIM1 protein into a P-TEFb (CDK9/cyclin T)inhibitor. *EMBO J.* 23, 2608–2619.

- Mignone, F., Gissi, C., Liuni, S., Pesole, G., 2002. Untranslated regions of mRNAs. *Genome Biol.* 3 (3) (reviews0004.1–0004.10).
- Mise, N., Goto, Y., Nakajima, N., Takagi, N., 1999. Molecular cloning of antisense transcripts of the mouse *Xist* gene. *Biochem. Biophys. Res. Commun.* 258, 537–541.
- Mishra, R.K., Eliceiri, G.L., 1997. Three small nucleolar RNAs that are involved in ribosomal RNA precursor processing. *Proc. Natl. Acad. Sci. USA* 94, 4972–4977.
- Mitchell, J.R., Cheng, J., Collins, K., 1999. A box H/ACA small nucleolar RNA-like domain at the human telomerase 3' end. *Mol. Cell Biol.* 19, 567–576.
- Miyata, T., Yasunaga, T., Nishida, T., 1980. Nucleotide sequence divergence and functional constraints in mRNA evolution. *Genetics* 77 (12), 7328–7332.
- Mochizuki, K., Fine, N.A., Fujisawa, T., Gorovsky, M.A., 2002. Analysis of a piwi-related gene implicates small RNAs in genome rearrangement in tetrahymena. *Cell* 110, 689–699.
- Møller-Jensen, J., Franch, T., Gerdes, K., 2001. Temporal translation control by metastable RNA structure. *J. Biol. Chem.* 276, 35707–35713.
- Montzka Wassarman, K., Storz, G., 2000. 6S RNA regulates *E. coli* RNA polymerase activity. *Cell* 101, 613–623.
- Moore, P.B., Steitz, T.A., 2002. The involvement of RNA in ribosome function. *Nature* 418, 229–235.
- Morey, C., Avner, P., 2004. Employment opportunities for non-coding RNAs. *FEBS Lett.* 567, 27–34.
- Morgenstern, B., 1999. DIALIGN2: improvement of the segment-to-segment approach to multiple sequence alignment. *Bioinformatics* 15, 211–218.
- Morrissey, J.P., Tollervey, D., 1995. Birth of the snoRNPs: the evolution of RNase MRP and the eukaryotic pre-rRNA-processing system. *Trends Biol. Sci.* 20, 78–82.
- Mosig, A., Sameith, K., Stadler, P.F., 2004. *fragrep*: efficient search for fragmented patterns in genomic sequences. Preprint, submitted for publication.
- Mourelatos, Z., Dostie, J., Paushkin, S., Sharma, A., Charroux, B., Abel, L., Rappaport, J., Mann, M., Dreyfuss, G., 2002. miRNPs: a novel class of ribonucleoproteins containing numerous microRNAs. *Genes Dev.* 16, 720–728.
- Myslinski, E., Krol, A., Carbon, P., 2004. Characterization of snRNA and snRNA-type genes in the pufferfish *Fugu rubripes*. *Gene* 330, 149–158.
- Nag, M.K., Thai, T.T., Ruff, E.A., Selvamurugan, N., Kunnumalaiyaan, M., Eliceiri, G.L., 1993. Genes for E1, E2, and E3 small nucleolar RNAs. *Proc. Natl. Acad. Sci. USA* 90, 9001–9005.
- Nagel, J.H.A., Gulyaev, A.P., Gerdes, K., Pleij, C.W.A., 1999. Metastable structures and refolding kinetics in *hok* mRNA of plasmid R1. *RNA* 5, 1408–1419.
- Nahvi, A., Barrick, J.E., Breaker, R.R., 2004. Coenzyme *b*<sub>12</sub> riboswitches are widespread genetic control elements in prokaryotes. *Nucleic Acids Res.* 32, 143–150.
- Nelson, P., Kiriakidou, M., Sharma, A., Maniataki, E., Mourelatos, Z., 2003. The microRNA world: small is mighty. *Trends Biochem. Sci.* 28, 534–540.
- Nilsen, T.W., 2001. Evolutionary origin of SL-addition *trans*-splicing: still an enigma. *Trends Genet.* 17, 678–680.
- Nilsen, T.W., 2003. The spliceosome: the most complex molecular machine in the cell? *Bioessays* 25, 1147–1149.
- Nitta, I., Kamada, Y., Noda, H., Ueda, T., Watanabe, K., 1998. Reconstitution of peptide bond formation with *Escherichia coli* 23S ribosomal RNA domains. *Science* 281, 666–669.
- Nudler, E., Mironov, A.S., 2004. The riboswitch control of bacterial metabolism. *Trends Biochem. Sci.* 29 (1), 11–17.
- O'Brien, C.A., Margelot, K., Wolin, S.L., 1993. Xenopus Ro ribonucleoproteins: members of an evolutionarily conserved class of cytoplasmic ribonucleoproteins. *Proc. Natl. Acad. Sci. USA* 90, 7250–7254.
- Oguchi, K., Tamura, K., Takahashi, H., 2004. Characterization of *Oryza sativa* telomerase reverse transcriptase and possible role of its phosphorylation in the control of telomerase activity. *Gene* 342, 57–66.
- Ohno, M., Mattaj, I., 1999. Meiosis: MeiRNA hits the spot. *Curr. Biol.* 28, R66–R69.
- Oleynikov, Y., Singer, R.H., 1998. RNA localization: different zipcodes, same postman? *Trends Cell Biol.* 8, 381–383.
- Olsen, G.J., Woese, C.R., 1993. Ribosomal RNA: a key to phylogeny. *FASEB J.* 7, 113–123.
- Omer, A., Lowe, T., Russel, A., Ebhardt, H., Eddy, S., Dennis, P., 2000. Homologs of small nucleolar RNAs in Archaea. *Science* 288, 517–522.
- Omoto, S., Ito, M., Tsutsumi, Y., Ichikawa, Y., Okuyama, H., Andi Brisibe, E., Saksena, N.K., Fuji, Y., 2004. HIV-1 nef suppression by virally encoded microRNA. *Retrovirology* 1, 44 (Epub).

- Otsuka, J., Sugaya, N., 2003. Advanced formulation of base pair changes in the stem regions of ribosomal RNAs its application to mitochondrial rRNAs: for resolving the phylogeny of animals. *J. Theor. Biol.* 222, 447–460.
- Pang, K.C., Stephen, S., Engström, P.G., Tajul-Arifin, K., Chen, W., Wahlestedt, C., Lenhard, B., Hayashizaki, Y., Mattick, J.S., 2005. RNAdb—comprehensive mammalian noncoding RNA database. *Nucleic Acids Res.* 33 (Database issue), D125–D130.
- Panopoulou, G., Hennig, S., Groth, D., Krause, A., Pousta, A.J., Herwig, R., Vingron, M., Lehrach, H., 2003. New evidence for genome-wide duplications at the origin of vertebrates using an amphioxus gene set and completed animal genomes. *Genome Res.* 13, 1056–1066.
- Pasquinelli, A.E., Reinhart, B.J., Slack, F., Martindale, M.Q., Kurodak, M.I., Maller, B., Hayward, D.C., Ball, E.E., Degnan, B., Müller, P., Spring, J., Srinivasan, A., Fishman, M., Finnerty, J., Corbo, J., Levine, M., Leahy, P., Davidson, E., Ruvkun, G., 2000. Conservation of the sequence and temporal expression of let-7 heterochronic regulatory RNA. *Nature* 408, 86–89.
- Pasquinelli, A.E., McCoy, A., Jiménez, E., Emili, S., Ruvkun, G., Martindale, M.Q., Baguñà, J., 2003. Expression of the 22 nucleotide let-7 heterochronic RNA throughout the metazoa: a role in life history evolution? *Evol. Dev.* 5, 372–378.
- Patel, A.A., Steitz, J.A., 2003. Splicing double: insights from the second spliceosome. *Nat. Rev. Mol. Cell Biol.* 4, 960–970.
- Paule, M.R., White, R.J., 2000. Survey and summary: transcription by RNA polymerases i and iii. *Nucleic Acids Res.* 28, 1283–1298.
- Pavesi, G., Mauri, G., Stefanini, M., Pesole, G., 2004. RNAProfile: an algorithm for finding conserved secondary structure motifs in unaligned RNA sequences. *Nucleic Acids Res.* 32, 3258–3269.
- Pedersen, J.S., Meyer, I.M., Forsberg, R., Hein, J., 2004a. An evolutionary model for protein-coding regions with conserved RNA structure. *Mol. Biol. Evol.* 21, 1913–1922.
- Pedersen, J.S., Meyer, I.M., Forsberg, R., Simmonds, P., Hein, J., 2004b. A comparative method for finding and folding RNA secondary structures within protein-coding regions. *Nucleic Acids Res.* 32, 4925–4936.
- Penny, D., Poole, A., 1999. The nature of the last universal common ancestor. *Curr. Opin. Genet. Dev.* 9, 672–677.
- Penny, G.D., Kay, G.F., Sheardown, S.A., Rastan, S., Brockdorff, N., 1996. The *Xist* gene is required in *cis* for X chromosome inactivation. *Nature* 379, 131–137.
- Pesole, G., Gissi, C., Grillo, G., Licciulli, F., Liuni, S., Saccone, C., 2000a. Analysis of oligonucleotide AUG start codon context in eukaryotic mRNAs. *Gene* 261, 85–91.
- Pesole, G., Liuni, S., D’Souza, M., 2000b. PatSearch: a pattern matcher software that finds functional elements in nucleotide and protein sequences and assesses their statistical significance. *Bioinformatics* 16 (5), 439–450.
- Pesole, G., Mignone, F., Gissi, C., Grillo, G., Licciulli, F., Liuni, S., 2001. Structural and functional features of eukaryotic mRNA untranslated regions. *Gene* 276, 73–81.
- Pesole, G., Liuni, S., Grillo, G., Licciulli, F., Mignone, F., Gissi, C., Saccone, C., 2002. UTRdb and UTRSite: specialized databases of sequences and functional elements of 5' and 3' untranslated regions of eukaryotic mRNAs. *Nucleic Acids Res.* 30, 335–340.
- Peterson, K., Eernisse, D.J., 2001. Animal phylogeny and the ancestry of bilaterians: inferences from morphology and 18S DNA gene sequences. *Evol. Dev.* 3, 170–205.
- Pfeffer, S., Zavolan, M., Grasser, F.A., Chien, M., Russo, J.J., Ju, J., John, B., Enright, A.J., Marks, D., Sander, C., Tuschl, T., 2004. Identification of virus-encoded microRNAs. *Science* 304, 734–736.
- Phillips, S.C., Turner, P.C., 1991. Sequence and expression of a mouse U7 snRNA type II pseudogene. *DNA Sequence* 1, 401–404.
- Pirotta, V., 2002. Trans-splicing in drosophila. *Bioessays* 24, 988–991.
- Pitulle, C., Garcia-Paris, M., Zamudio, K.R., Pace, N.R., 1998. Comparative structural analysis of vertebrate ribonuclease P RNA. *Nucleic Acids Res.* 26, 3333–3339.
- Pokrywka, N.J., Stephenson, E.C., 1991. Microtubules mediate the localization of bicoid RNA during *Drosophila* oogenesis. *Development* 113, 55–66.
- Poole, A., Penny, D., Sjäberg, B.-M., 2000. Methyl-RNA: an evolutionary bridge between RNA and DNA? *Chem. Biol.* 7, R207–R216.
- Potter, S.S., Branford, W.W., 1998. Evolutionary conservation and tissue-specific processing of *Hoxa 11* antisense transcripts. *Mamm. Genome* 9, 799–806.
- Precott, E.M., Proudfoot, N.J., 2002. Transcriptional collision between convergent genes in budding yeast. *Proc. Natl. Acad. Sci. USA* 99, 8796–8801.

- Prohaska, S.J., Fried, C., Flamm, C., Wagner, G.P., Stadler, P.F., 2004. Surveying phylogenetic footprints in large gene clusters: applications to Hox cluster duplications. *Mol. Phyl. Evol.* 31, 581–604.
- Putzer, H., Gendron, N., Grunberg-Manago, M., 1992. Co-ordinate expression of the two threonyl-tRNA synthetase genes in *Bacillus subtilis*: control by transcriptional antitermination involving a conserved regulatory sequence. *EMBO J.* 11, 3117–3127.
- Ramakrishnan, V., Moore, P.B., 2001. Atomic structures at last: the ribosome in 2000. *Curr. Opin. Struct. Biol.* 11, 144–154.
- Regalado, M., Rosenblad, M.A., Samuelson, T., 2002. Prediction of signal recognition particle RNA genes. *Nucleic Acids Res.* 30, 3368–3377.
- Rehmsmeier, M., Steffen, P., Höchsmann, M., Giegerich, R., 2004. Fast and effective prediction of microRNA/target duplexes. *RNA* 10, 1507–1517.
- Reinhart, F.J., Slack, B.J., Basson, M., Pasquinelli, A.E., Bettinger, J.C., Rougvie, A.E., Horwitz, H.R., Ruvkun, G., 2000. The 21-nucleotide RNA *let-7* regulates developmental timing in *Caenorhabditis elegans*. *Nature* 403, 901–906.
- Rivas, E., Eddy, S.R., 2000. Secondary structure alone is generally not statistically significant for the detection of noncoding RNAs. *Bioinform.* 16, 583–605.
- Rivas, E., Eddy, S.R., 2001. Noncoding RNA gene detection using comparative sequence analysis. *BMC Bioinform.* 2, 8.
- Rivas, E., Klein, R.J., Jones, T.A., Eddy, S.R., 2001. Computational identification of noncoding RNAs in *E. coli* by comparative genomics. *Curr. Biol.* 11, 1369–1373.
- Rodin, S., Ohno, S., Rodin, A., 1993. Transfer RNAs with complementary anticodons: could they reflect early evolution of discriminative genetic code adaptors? *Proc. Natl. Acad. Sci. USA* 90, 4723–4727.
- Rodionov, D.A., Vitreschak, A.G., Mironov, A.A., Gelfand, M.S., 2004. Comparative genomics of the methionine metabolism in Gram-positive bacteria: a variety of regulatory systems. *Nucleic Acids Res.* 32, 3340–3353.
- Rodriguez, A., Griffiths-Jones, S., Ashurst, J.L., Bradley, A., 2004. Identification of mammalian microRNA host genes and transcription units. *Genome Res.* 14, 1902–1910.
- Rooney, A.P., 2004. Mechanisms underlying the evolution and maintenance of functionally heterogeneous 18S rRNA genes in apicomplexans. *Mol. Biol. Evol.* 21, 1704–1711.
- Rosenblad, M.A., Samuelsson, T., 2004. Identification of chloroplast signal recognition particle RNA genes. *Plant Cell Physiol.* 45, 1633–1639.
- Rosenblad, M.A., Gorodkin, J., Knudsen, B., Zwieb, C., 2003. SRPDB: signal recognition particle database. *Nucleic Acids Res.* 31, 363–364.
- Rosenblad, M.A., Zwieb, C., Samuelson, T., 2004. Identification and comparative analysis of components from the signal recognition particle in protozoa and fungi. *BMC Genomics* 5 (5).
- Rueckert, R.R., 1996. Picornaviridae: the viruses and their replication. In: Fields, N., Knipe, D., Howley, P. (Eds.), *Virology*, vol. 1. third ed. Lippincott-Raven Publishers, Philadelphia, New York, pp. 609–654.
- Russell, A.G., Schnare, M.N., Gray, M.W., 2004. Pseudouridine-guide RNAs and other Cbf5p-associated RNAs in *Euglena gracilis*. *RNA* 10, 1034–1046.
- Saitou, N., Nei, M., 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol. Biol. Evol.* 4, 406–425.
- Samarsky, D.A., Fournier, M.J., 1999. A comprehensive database for the small nucleolar RNAs from *Saccharomyces cerevisiae*. *Nucleic Acids Res.* 27, 161–164.
- Samarsky, D.A., Schneider, G.S., Fournier, M.J., 1996. An essential domain in *Saccharomyces cerevisiae* U14 snoRNA is absent in vertebrates, but conserved in other yeasts. *Nucleic Acids Res.* 24, 2059–2066.
- Sankoff, D., 1985. Simultaneous solution of the RNA folding, alignment, and proto-sequence problems. *SIAM J. Appl. Math.* 45, 810–825.
- Savill, N.J., Hoyle, D.C., Higgs, P.G., 2001. RNA sequence evolution with secondary structure constraints: comparison of substitution rate models using maximum-likelihood methods. *Genetics* 157, 399–411.
- Scharl, E.C., Steitz, J.A., 1996. Length suppression in histone messenger RNA 3'-end maturation: processing defects of insertion mutant pre-messenger RNAs can be compensated by insertions into the U7 small nuclear RNA. *Proc. Natl. Acad. Sci. USA* 93, 14659–14664.
- Schattner, P., 2002. Searching for RNA genes using base-composition statistics. *Nucleic Acids Res.* 30 (9), 2076–2082.
- Schlötterer, C., Tautz, D., 1994. Chromosomal homogeneity of drosophila ribosomal DNA arrays suggests intrachromosomal exchanges drive concerted evolution. *Curr. Biol.* 4, 777–783.

- Schöninger, M., von Haeseler, A., 1999. Towards assigning helical regions in alignments of ribosomal RNA and testing the appropriateness of evolutionary models. *J. Mol. Evol.* 49, 691–698.
- Schramm, L., Hernandez, N., 2002. Recruitment of RNA polymerase III to its target promoters. *Genes Dev.* 16, 2593–2620.
- Schultes, E.A., Bartel, D.P., 2000. One sequence, two ribozymes: implications for the emergence of new ribozyme folds. *Science* 289, 448–452.
- Schultes, E.A., Hrabr, P.T., LaBean, T.H., 1999. Estimating the contributions of selection and self-organization in RNA secondary structure. *J. Mol. Evol.* 49, 76–83.
- Schümperli, D., Pillai, R.S., 2004. The special Sm core structure of the U7 snRNP: far-reaching significance of a small nuclear ribonucleoprotein. *Cell. Mol. Life Sci.* 61, 2560–2570.
- Schuster, P., Fontana, W., Stadler, P.F., Hofacker, I.L., 1994. From sequences to shapes and back: a case study in RNA secondary structures. *Proc. R. Soc. London B* 255, 279–284.
- Seitz, H., Royo, H., Lin, S.-P., Youngson, N., Ferguson-Smith, A.C., Cavaillé, J., 2004. Imprinted small RNA genes. *Biol. Chem.* 385, 905–911.
- Selvamurugan, N., Eliceiri, G.L., 1995. The gene for human E2 small nucleolar RNA resides in an intron of a laminin-binding protein gene. *Genomics* 30, 400–401.
- Shabalina, S.A., Ogurtsov, A.Y., Rogozin, I.B., Koonin, E.V., Lipman, D.J., 2004. Comparative analysis of orthologous eukaryotic mRNAs: potential hidden functional signals. *Nucleic Acids Res.* 32, 1774–1782.
- Shapiro, B.A., Zhang, K., 1990. Comparing multiple RNA secondary structures using tree comparisons. *CABIOS* 6, 309–318.
- Sharkady, S.M., Williams, K.P., 2004. A third lineage with two-piece tmRNA. *Nucleic Acids Res.* 32, 1–8.
- Shendure, J., Church, G.M., 2002. Computational discovery of senseantisense transcription in the human and mouse genome. *Genome Biol.* 3, 1–14.
- Siebert, S., Backofen, R., 2003. MARNA: a server for multiple alignment of RNAs. In: Mewes, H.-W., Heun, V., Frishman, D., Kramer, S., (Eds.), proceedings of the German Conference on Bioinformatics. GCB 2003. vol. 1, München, D. Belleville Verlag, Michael Farin, pp. 135–140.
- Silverman, S.K., 2003. Rube goldberg goes (ribo)nuclear? Molecular switches and sensors made from RNA. *RNA* 9, 377–383.
- Simpson, L., Thiemann, O.H., Savill, N.J., Alfonzo, J.D., Maslov, D.A., 2000. Evolution of RNA editing in trypanosome mitochondria. *Proc. Natl. Acad. Sci. USA* 97, 6986–6993.
- Soldati, D., Schümperli, D., 1988. Structural and functional characterization of mouse U7 small nuclear RNA active in 3' processing of histone pre-mRNA. *Mol. Cell Biol.* 8, 1518–1524.
- Soukup, G.A., Breaker, R.R., 1999. Engineering precision RNA molecular switches. *Proc. Natl. Acad. Sci. USA* 96, 3584–3589.
- Sprinzl, M., Vassilenko, K.S., 2005. Compilation of tRNA sequences and sequences of tRNA genes. *Nucleic Acids Res.* 33 (Database issue), 139–140.
- Sprinzl, M., Steegborn, C., Hübel, F., Steinberg, S., 1996. Compilation of tRNA sequences and sequences of tRNA genes. *Nucleic Acids Res.* 24, 68–72.
- Sprinzl, M., Horn, C., Brown, M., Ioudovitch, A., Steinberg, S., 1998. Compilation of tRNA sequences and sequences of tRNA genes. *Nucleic Acids Res.* 26, 148–153.
- Steitz, T.A., Moore, P.B., 2003. RNA, the first macromolecular catalyst: the ribosome is a ribozyme. *Trends Biochem. Sci.* 28, 411–418.
- Storz, G., Opdyke, J.A., Zhang, A., 2004. Controlling mRNA stability and translation with small noncoding RNAs. *Curr. Opin. Microbiol.* 7, 140–144.
- Stuart, K., Allen, T.E., Heidmann, S., Seiwert, S.D., 1997. RNA editing in kinetoplastid protozoa. *Microbiol. Mol. Biol. Rev.* 61, 105–120.
- Sudarsan, N., Barrick, J.E., Breaker, R.R., 2003. Metabolite-binding RNA domains are present in the genes of eukaryotes. *RNA* 9, 644–647.
- Sullenger, B.A., 2004. Riboswitches—to kill or save the messenger. *N. Engl. J. Med.* 351, 2759–2760.
- Sumiyama, K., Irvine, S.Q., Ruddle, F.H., 2003. The role of gene duplication in the evolution and function of the vertebrate Dlx/distal-less bigene clusters. *J. Struct. Funct. Genomics* 3, 151–159.
- Sun, Y., Koo, S., White, N., Peralta, E., Esau, C., Dean, N.M., Perera, R.J., 2004. Development of a micro-array to detect human and mouse microRNAs and characterization of expression in human organs. *Nucleic Acids Res.* 32 (doi:10.1093/nar/gnh186).
- Suzuki, M., Hayashizaki, Y., 2004. Mouse-centric comparative transcriptomics of protein coding and non-coding RNAs. *BioEssays* 26, 833–843.
- Szymański, M., Barciszewska, M., Barciszewski, J., Erdmann, V., 2000. 5S ribosomal RNA database Y2K. *Nucleic Acids Res.* 28, 166–167.

- Szymański, M., Barciszewska, M.Z., Żywicki, M., Barciszewski, J., 2003. Noncoding RNA transcripts. *J. Appl. Genet.* 44, 1–19.
- Talla, E., Anthouard, V., Bouchier, C., Frangeul, L., Dujon, B., 2005. The complete mitochondrial genome of the yeast *Kluyveromyces thermotolerans*. *FEBS Lett.* 579, 30–40.
- Tang, T.-H., Bachellerie, J.-P., Rozhdestvensky, T., Bortolin, M.-L., Huber, H., Drungowski, M., Elge, T., Brosius, J., Hüttnerhofer, A., 2002. Identification of 86 candidates for small non-messenger RNAs from the archaeon *Archaeoglobus fulgidus*. *Proc. Natl. Acad. Sci. USA* 99, 7536–7541.
- Tanzer, A., Stadler, P.F., 2004. Molecular evolution of a microRNA cluster. *J. Mol. Biol.* 339, 327–335.
- Tanzer, A., Amemiya, C.T., Kim, C.-B., Stadler, P.F., 2005. Evolution of microRNAs located within Hox gene clusters. *J. Exp. Zool.: Mol. Dev. Evol.* 304B, 75–85.
- Tarn, W.Y., Yario, T.A., Steitz, J.A., 1995. U12 snRNAs in vertebrates: evolutionary conservation of 5' sequences implicated in splicing of pre-mRNAs containing a minor class of introns. *RNA* 1, 644–656.
- Telford, M.J., Holland, P.W.H., 1997. Evolution of 28S ribosomal DNA in chaetognaths: duplicate genes and molecular phylogeny. *J. Mol. Evol.* 44, 135–144.
- Terns, M.P., Terns, R.M., 2002. Small nucleolar RNAs: versatile transacting molecules of ancient evolutionary origin. *Gene Exp.* 10, 17–39.
- Teunissen, S.W., Kruithof, M.J., Farris, A.D., Harley, J.B., Venrooij, W.J., Pruijn, G.J., 2000. Conserved features of Y RNAs: a comparison of experimentally derived secondary structures. *Nucleic Acids Res.* 28, 610–619.
- Thompson, J.D., Higgins, D.G., Gibson, T.J., 1994. CLUSTALW: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* 22, 4673–4680.
- Thompson, M., Haeusler, R.A., Good, P.D., Engelke, D.R., 2003. Nucleolar clustering of dispersed tRNA genes. *Science* 302, 1399–1401.
- Thurner, C., Witwer, C., Hofacker, I., Stadler, P.F., 2004. Conserved RNA secondary structures in Flaviviridae genomes. *J. Gen. Virol.* 85, 1113–1124.
- Tran, E., Brown, J., Maxwell, S.E., 2004. Evolutionary origins of the RNA-guided nucleotide-modification complexes: from the primitive translation apparatus? *Trends Biochem. Sci.* 29, 343–350.
- Tschudi, C., Ullu, E., 2002. Unconventional rules of small nuclear RNA transcription and cap modification in trypanosomatids. *Gene Exp.* 10, 3–16.
- Tuplin, A., Wood, J., Evans, D.J., Patel, A.H., Simmonds, P., 2002. Thermodynamic and phylogenetic prediction of RNA secondary structures in the coding region of hepatitis C virus. *RNA* 8, 824–841.
- Tuplin, A., Evans, D.J., Simmonds, P., 2004. Detailed mapping of RNA secondary structures in core and NS5B-encoding region sequence of hepatitis C virus by RNase cleavage and novel bioinformatic prediction methods. *J. Gen. Virol.* 85, 3037–3047.
- Turner, I.A., Norman, C.M., Churcher, M.J., N.A. J., 2004. Roles of the U5 snRNP in spliceosome dynamics and catalysis. *Biochem. Soc. Trans.* 32, 928–931.
- Tycowski, K.T., Steitz, J.A., 2001. Non-coding snoRNA host genes in *Drosophila*: expression strategies for modification guide snoRNAs. *Eur. J. Cell. Biol.* 80, 119–125.
- Tycowski, K.T., Aab, A., Steitz, J.A., 2004. Guide RNAs with 5' caps and novel box C/D snoRNA-like domains for modification of snRNAs in metazoa. *Curr. Biol.* 14, 1985–1995.
- Uliel, S., Liang, X.-h., Unger, R., Michaeli, S., 2004. Small nucleolar RNAs that guide modification in trypanosomatids: repertoire, targets, genome organization, and unique functions. *Int. J. Parasitol.* 34, 445–454.
- Unrau, P.J., Bartel, D.P., 1998. RNA-catalysed nucleotide synthesis. *Nature* 395, 260–263.
- Ushida, C., Yoshida, A., Miyakawa, Y., Ara, Y., Muto, A., 2003. Distribution of the MCS4 RNA genes in mycoplasmas belonging to the *Mycoplasma mycoides* cluster. *Gene* 314, 149–155.
- Valadkhan, S., Manley, J.L., 2001. Splicing-related catalysis by protein-free snRNAs. *Nature* 413, 701–707.
- Valadkhan, S., Manley, J.L., 2003. Characterization of the catalytic activity of U2 and U6 snRNAs. *RNA* 9, 892–904.
- Van de Peer, Y., Baldauf, S.L., Doolittle, W.F., Meyer, A., 2000a. An updated and comprehensive rRNA phylogeny of (crown) eukaryotes based on rate-calibrated evolutionary distances. *J. Mol. Evol.* 51, 565–576.
- Van de Peer, Y., De Rijk, P., Wuylts, J., Winkelmann, T., DeWachter, R., 2000b. The european small subunit ribosomal RNA database. *Nucleic Acids Res.* 28, 175–176.
- Van Horn, D.J., Eisenberg, D., O'Brien, C.A., Wolin, S.L., 1995. *Caenorhabditis elegans* embryos contain only one major species of Ro RNP. *RNA* 1, 293–303.

- van Zon, A., Mossink, M., Schoester, M., Scheffer, G., Scheper, R., Sonneveld, P., Wiemer, E., 2001. Multiple human vault RNAs. Expression and association with the vault complex. *J. Biol. Chem.* 276, 37715–37721.
- Vasu, S.K., Rome, L.H., 1995. *Dictyostelium* vaults: disruption of the major proteins reveals growth and morphological defects and uncovers a new associated protein. *J. Biol. Chem.* 270, 16588–16594.
- Vella, M.C., Reinert, K., Slack, F.J., 2004. Architecture of a validated MicroRNA: target interaction. *Chem. Biol.* 11, 1619–1623.
- Vitali, P., Royo, H., Seitz, H., Bachellerie, J.-P., Hüttenhofer, A., Cavallé, J., 2003. Identification of 13 novel human modification guide RNAs. *Nucleic Acids Res.* 31, 6543–6551.
- Vitreschak, A.G., Rodionov, D.A., Mironov, A.A., Gelfand, M.S., 2003. Regulation of the vitamine B<sub>12</sub> metabolism and transport in bacteria by a conserved RNA structural element. *RNA* 9, 1084–1097.
- Vitreschak, A.G., Rodionov, D.A., Mironov, A.A., Gelfand, M.S., 2004. Riboswitches: the oldest mechanism for the regulation of gene expression? *Trends Genet.* 20 (1), 44–50.
- Vogel, J., Bartels, V., Tang, T.H., Churakov, G., Slagter-Jäger, J.G., Hüttenhofer, A., Wagner, G.H.E., 2003. RNomics in *Escherichia coli* detects new sRNA species and indicates parallel transcriptional output in bacteria. *Nucleic Acids Res.* 31, 6435–6443.
- Wagner, E.G.H., Flärdh, K., 2002. Antisense RNAs everywhere? *Trends Genet.* 18, 223–226.
- Washietl, S., Hofacker, I.L., 2004. Consensus folding of aligned sequences as a new measure for the detection of functional RNAs by comparative genomics. *J. Mol. Biol.* 342, 19–30.
- Washietl, S., Hofacker, I.L., Stadler, P.F., 2005. Fast and reliable detection of noncoding RNAs. *Proc. Natl. Acad. Sci.* 102, 2454–2459.
- Wassarman, D.A., Steitz, J.A., 1991. Structural analyses of the 7SK ribonucleoprotein (RNP), the most abundant human small RNP of unknown function. *Mol. Cell. Biol.* 11, 3432–3445.
- Wassarman, K., Repoila, F., Rosenow, C., Storz, G., Gottesman, S., 2001. Identification of novel small RNAs using comparative genomics and microarrays. *Genes Dev.* 15, 1637–1651.
- Weber, M.J., 2005. New human and mouse microRNA genes found by homology search. *FEBS J.* 272, 59–73.
- Weiner, A.M., Denison, R.A., 1983. Either gene amplification or gene conversion may maintain the homogeneity of the multigene family encoding human U1 small nuclear RNA. *Cold Spring Harber Symp. Quant. Biol.* 47, 1141–1149.
- Weinstein, L.B., Steitz, J.A., 1999. Guided tours: from precursor snoRNA to functional snoRNP. *Curr. Opin. Cell Biol.* 11, 378–384.
- Werner, A., Preston-Fayers, K., Dehmelt, L., Nalbant, P., 2002. Regulation of the NPT gene by a naturally occurring antisense transcript. *Cell Biochem. Biophys.* 36, 241–252.
- Westhof, E., Massire, C., 2004. Evolution of RNA architecture. *Science* 306, 62–63.
- White, R.J., 1998. RNA Polymerase III Transcription. Springer, New York, NY.
- Wilkie, G.S., Dickson, K.S., Gray, N.G., 2003. Regulation of mRNA translation by 5'- and 3'-UTR-binding factors. *Trends Biochem. Sci.* 28, 182–188.
- Williams, K.P., 2002. Descent of a split DNA. *Nucleic Acids Res.* 30, 2025–2030.
- Winkler, W.C., Breaker, R.R., 2003. Genetic control by metabolite-binding riboswitches. *Chembiochem.* 4 (10), 1024–1032.
- Witwer, C., Rauscher, S., Hofacker, I., Stadler, P., 2001. Conserved RNA secondary structures in picornaviridae genomes. *Nucleic Acids Res.* 29, 5079–5089.
- Wolffe, A.P., 1994. The role of transcription factors, chromatin structure and DNA replication in 5S RNA gene regulation. *J. Cell Sci.* 107, 2055–2063.
- Wood, V., Gwilliam, R., Rajandream, M.A., et al., 2002. (132 co-authors). The genome sequence of *Schizosaccharomyces pombe*. *Nature* 415, 871–880.
- Wu, C.-H.H., Gall, J.G., 1993. U7 small nuclear RNA in C snurposomes of the Xenopus germinal vesicle. *Proc. Natl. Acad. Sci. USA* 90, 6257–6259.
- Wuyts, J., De Rijk, P., Van de Peer, Y., Winkelmann, T., De Wachter, R., 2001. The european large subunit ribosomal RNA database. *Nucleic Acids Res.* 29, 175–177.
- Yao, M.C., Fuller, P., Xi, X., 2003. Programmed DNA deletion as an RNA-guided system of genome defense. *Science* 300, 1517–1518.
- Ye, A.J., Romero, D.P., 2002. Phylogenetic relationships amongst tetrahymenine ciliates inferred by a comparison of telomerase RNAs. *Int. J. Syst. Evol. Microbiol.* 52, 2297–2302.
- Yekta, S., Shih, I.-H., Bartel, D.P., 2004. MircoRNA-directed cleavage of *HoxB8* mRNA. *Science* 304, 594–596.

- Yelin, R., Dahary, D., Sorek, R., Levanon, E.Y., Goldstein, O., Shoshan, A., Diber, A., Biton, S., Tamir, Y., Khosravi, R., Nemzer, S., Pinner, E., Walach, S., Bernstein, J., Savitsky, K., Rotman, G., 2003. Widespread occurrence of antisense transcription in the human genome. *Nat. Biotechnol.* 21, 379–386.
- Yik, J.H., Chen, R., Nishimura, R., Jennings, J.L., Link, A.J., Zhou, Q., 2003. Inhibition of P-TEFb (CDK9/Cyclin T) kinase and RNA polymerase II transcription by the coordinated actions of HEXIM1 and 7SK snRNA. *Mol. Cell* 12, 971–982.
- Ying, S.-Y., Lin, S.-L., 2004. Intron-derived microRNAs—fine tuning of gene functions. *Gene* 342, 25–28.
- Ying, S.-Y., Lin, S.-L., 2005. Intronic microRNAs. *Biochem. Biophys. Res. Comm.* 326, 515–520.
- Yu, Y.-T., Tarn, W.-Y., Yario, T.A., Steitz, J.A., 1996. More Sm snRNAs from vertebrate cells. *Exp. Cell Res.* 229, 276–281.
- Zappulla, D.C., Cech, T.R., 2004. Yeast telomerase RNA: a flexible scaffold for protein subunits. *Proc. Natl. Acad. Sci. USA* 101, 10024–10029.
- Zhang, Z., Gerstein, M., 2003. Of mice and men: phylogenetic footprinting aids the discovery of regulatory elements. *J. Biol.* 2 (11, 4pp.).
- Zimmerly, S., Hausner, G., Wu, X.-C., 2001. Phylogenetic relationships among group II intron ORFs. *Nucleic Acids Res.* 29, 1238–1250.
- Zwieb, C., Eichler, J., 2002. Getting on target: the archaeal signal recognition particle. *Archaea* 1, 27–34.
- Zwieb, C., Wower, J., 2000. tmRDB (tmRNA database). *Nucleic Acids Res.* 28, 169–170.
- Zwieb, C., van Nues, R.W., Rosenblad, M.A., Brown, J.D., Samuelson, T., 2005. A nomenclature for all signal recognition particle RNAs. *RNA* 11, 7–13.